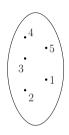
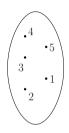
Learning Distributions using Lie Groups

E. Mehmet Kıral

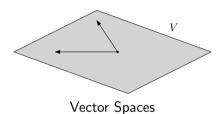
August 2025, IJCAI tutorial AI meets Algebra¹

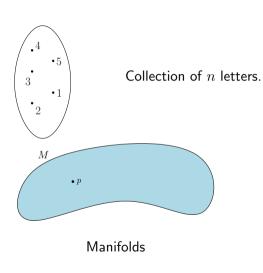


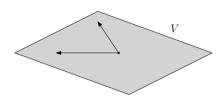
Collection of n letters.



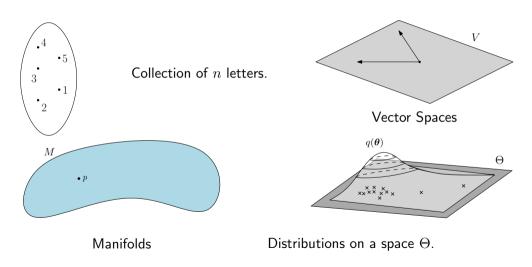
Collection of n letters.



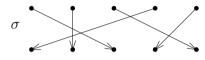


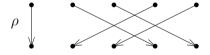


Vector Spaces

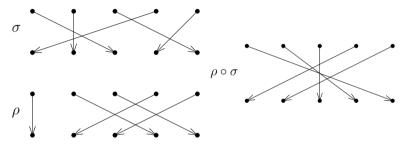


Permutations on n letters. This preserves the distinctness of elements of the set. Given two permutations you can apply one after the other.

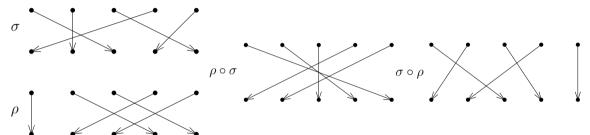




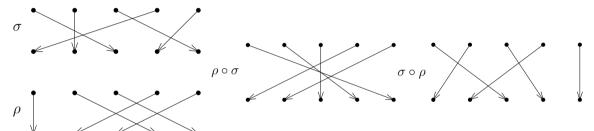
Permutations on n letters. This preserves the distinctness of elements of the set. Given two permutations you can apply one after the other.



Permutations on n letters. This preserves the distinctness of elements of the set. Given two permutations you can apply one after the other.



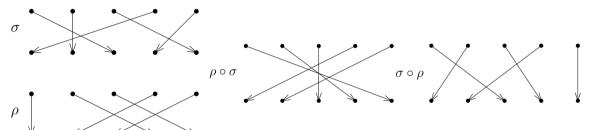
Permutations on n letters. This preserves the distinctness of elements of the set. Given two permutations you can apply one after the other.



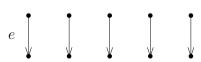
The identity permutation:



Permutations on n letters. This preserves the distinctness of elements of the set. Given two permutations you can apply one after the other.



The identity permutation:



The resulting object is the symmetric group on n letters, denoted by S_n or Sym_n .

Invertible linear maps also preserve the linear structure of vectors. It is a group.

$$\mathrm{GL}(V) = \left\{ A: V \to V \middle| egin{aligned} A(\mathbf{v} + \mathbf{w}) &= A\mathbf{v} + A\mathbf{w}, \\ A(\alpha\mathbf{v}) &= \alpha A\mathbf{v} \end{aligned}, A^{-1} \text{ exists} \right\}.$$

Also denoted by $GL(n,\mathbb{R})$ if V is an n-dimensional real vector space, i.e. $V \cong \mathbb{R}^n$.

Invertible linear maps also preserve the linear structure of vectors. It is a group.

$$\mathrm{GL}(V) = \left\{ A : V \to V \middle| egin{aligned} A(\mathbf{v} + \mathbf{w}) &= A\mathbf{v} + A\mathbf{w}, \\ A(\alpha\mathbf{v}) &= \alpha A\mathbf{v} \end{aligned}, A^{-1} \text{ exists} \right\}.$$

Also denoted by $\mathrm{GL}(n,\mathbb{R})$ if V is an n-dimensional real vector space, i.e. $V\cong\mathbb{R}^n$.

If preserve extra structure such as inner product $\langle \mathbf{v}, \mathbf{w} \rangle = \mathbf{v}^{\top} \mathbf{w}$ then obtain a subgroup

$$SO(n, \mathbb{R}) = \{ A \in GL(n, \mathbb{R}) | \langle A\mathbf{v}, A\mathbf{w} \rangle = \langle \mathbf{v}, \mathbf{w} \rangle, \det(A) > 0 \}.$$

Invertible linear maps also preserve the linear structure of vectors. It is a group.

$$\mathrm{GL}(V) = \left\{ A: V \to V \middle| \begin{array}{l} A(\mathbf{v} + \mathbf{w}) = A\mathbf{v} + A\mathbf{w}, \\ A(\alpha\mathbf{v}) = \alpha A\mathbf{v} \end{array}, A^{-1} \text{ exists} \right\}.$$

Also denoted by $\mathrm{GL}(n,\mathbb{R})$ if V is an n-dimensional real vector space, i.e. $V\cong\mathbb{R}^n$.

If preserve extra structure such as inner product $\langle \mathbf{v}, \mathbf{w} \rangle = \mathbf{v}^{\top} \mathbf{w}$ then obtain a subgroup

$$SO(n, \mathbb{R}) = \{ A \in GL(n, \mathbb{R}) | \langle A\mathbf{v}, A\mathbf{w} \rangle = \langle \mathbf{v}, \mathbf{w} \rangle, \det(A) > 0 \}.$$

The upper triangular matrices are also closed under multiplication and inversion

$$B = \begin{pmatrix} * & * & * & \cdots & * \\ 0 & * & * & & * \\ 0 & 0 & \ddots & & * \\ 0 & 0 & \cdots & 0 & * \end{pmatrix}$$

Invertible linear maps also preserve the linear structure of vectors. It is a group.

$$\mathrm{GL}(V) = \left\{ A: V \to V \middle| \begin{array}{l} A(\mathbf{v} + \mathbf{w}) = A\mathbf{v} + A\mathbf{w}, \\ A(\alpha\mathbf{v}) = \alpha A\mathbf{v} \end{array}, A^{-1} \text{ exists} \right\}.$$

Also denoted by $GL(n,\mathbb{R})$ if V is an n-dimensional real vector space, i.e. $V \cong \mathbb{R}^n$.

If preserve extra structure such as inner product $\langle \mathbf{v}, \mathbf{w} \rangle = \mathbf{v}^{\top} \mathbf{w}$ then obtain a subgroup

$$SO(n, \mathbb{R}) = \{ A \in GL(n, \mathbb{R}) | \langle A\mathbf{v}, A\mathbf{w} \rangle = \langle \mathbf{v}, \mathbf{w} \rangle, \det(A) > 0 \}.$$

The upper triangular matrices are also closed under multiplication and inversion

$$B = \begin{pmatrix} * & * & * & \cdots & * \\ 0 & * & * & & * \\ 0 & 0 & \ddots & & * \\ 0 & 0 & \cdots & 0 & * \end{pmatrix} \qquad \begin{array}{c} \text{What does this set of} \\ \text{matrices preserve?} \end{array}$$



Groups

A group is a set G together with a binary operation

$$*: G \times G \longrightarrow G$$

 $(g,h) \longmapsto g * h$

satisfying

- (associativity) (g*h)*k = g*(h*k) for every $g,h,k \in G$.
- ② (existence of identity) There is an $e \in G$ such that e * g = g and g * e = g for every $g \in G$.
- **3** (existence of inverses) For every $g \in G$ there exists an $h \in G$ such that g * h = h * g = e.

Groups

A group is a set G together with a binary operation

$$*: G \times G \longrightarrow G$$

 $(g,h) \longmapsto g * h$

satisfying

- (associativity) (g*h)*k = g*(h*k) for every $g,h,k \in G$.
- ② (existence of identity) There is an $e \in G$ such that e * g = g and g * e = g for every $g \in G$.
- **3** (existence of inverses) For every $g \in G$ there exists an $h \in G$ such that g * h = h * g = e.

Additive symbol '*' = '+' for Abelian groups: a + b = b + a.

Multiplicative notation mostly omits the * symbol, and $ab \neq ba$ can happen (think matrices).



Groups

A group is a set G together with a binary operation

$$*: G \times G \longrightarrow G$$

 $(g,h) \longmapsto g * h$

satisfying

- (associativity) (g*h)*k = g*(h*k) for every $g,h,k \in G$.
- ② (existence of identity) There is an $e \in G$ such that e * g = g and g * e = g for every $g \in G$.
- **③** (existence of inverses) For every $g \in G$ there exists an $h \in G$ such that g*h = h*g = e.

Additive symbol '*' = '+' for Abelian groups: a + b = b + a.

Multiplicative notation mostly omits the * symbol, and $ab \neq ba$ can happen (think matrices).

The closure property of the group operation is simply that $a * b \in G$ if $a, b \in G$.

 \bullet (V,+). Any vector space V with vector addition as the binary operation.

- \bullet (V,+). Any vector space V with vector addition as the binary operation.
- $(\mathbb{R}_{>0},\times)$. Positive real numbers $\mathbb{R}_{>0}$ under multiplication.

- (V, +). Any vector space V with vector addition as the binary operation.
- $(\mathbb{R}_{>0}, \times)$. Positive real numbers $\mathbb{R}_{>0}$ under multiplication.
- **3** $GL(n,\mathbb{R})$. Invertible $n \times n$ matrices, under matrix multiplication.

- \bullet (V,+). Any vector space V with vector addition as the binary operation.
- $(\mathbb{R}_{>0},\times)$. Positive real numbers $\mathbb{R}_{>0}$ under multiplication.
- **③** $\mathrm{GL}(n,\mathbb{R})$. Invertible $n \times n$ matrices, under matrix multiplication.
- **③** Aff $(n, \mathbb{R}) \cong \operatorname{GL}(n, \mathbb{R}) \ltimes \mathbb{R}^n$. The Affine group, consisting of pairs (A, \mathbf{b}) where the group multiplication is $(A_1, \mathbf{b}_1)(A_2, \mathbf{b}_2) = (A_1A_2, A\mathbf{b}_2 + \mathbf{b}_1)$.

- \bullet (V,+). Any vector space V with vector addition as the binary operation.
- $(\mathbb{R}_{>0},\times)$. Positive real numbers $\mathbb{R}_{>0}$ under multiplication.
- **③** $\mathrm{GL}(n,\mathbb{R})$. Invertible $n \times n$ matrices, under matrix multiplication.
- Aff $(n, \mathbb{R}) \cong \operatorname{GL}(n, \mathbb{R}) \ltimes \mathbb{R}^n$. The Affine group, consisting of pairs (A, \mathbf{b}) where the group multiplication is $(A_1, \mathbf{b}_1)(A_2, \mathbf{b}_2) = (A_1A_2, A\mathbf{b}_2 + \mathbf{b}_1)$. Another way to realize this group is as a subgroup of $(n+1) \times (n+1)$ matrices.

$$\begin{pmatrix} A_1 & \mathbf{b}_1 \\ \mathbf{0}^\top & 1 \end{pmatrix} \begin{pmatrix} A_2 & \mathbf{b}_2 \\ \mathbf{0}^\top & 1 \end{pmatrix} = \begin{pmatrix} A_1 A_2 & \mathbf{b}_1 + A_1 \mathbf{b}_2 \\ \mathbf{0}^\top & 1 \end{pmatrix} \text{ and } \begin{pmatrix} A & \mathbf{b} \\ \mathbf{0}^\top & 1 \end{pmatrix}^{-1} = \begin{pmatrix} A^{-1} & -A^{-1} \mathbf{b} \\ \mathbf{0}^\top & 1 \end{pmatrix}$$

- \bullet (V,+). Any vector space V with vector addition as the binary operation.
- $(\mathbb{R}_{>0},\times)$. Positive real numbers $\mathbb{R}_{>0}$ under multiplication.
- **③** $\mathrm{GL}(n,\mathbb{R})$. Invertible $n \times n$ matrices, under matrix multiplication.
- Aff $(n, \mathbb{R}) \cong \operatorname{GL}(n, \mathbb{R}) \ltimes \mathbb{R}^n$. The Affine group, consisting of pairs (A, \mathbf{b}) where the group multiplication is $(A_1, \mathbf{b}_1)(A_2, \mathbf{b}_2) = (A_1A_2, A\mathbf{b}_2 + \mathbf{b}_1)$. Another way to realize this group is as a subgroup of $(n+1) \times (n+1)$ matrices.

$$\begin{pmatrix} A_1 & \mathbf{b}_1 \\ \mathbf{0}^\top & 1 \end{pmatrix} \begin{pmatrix} A_2 & \mathbf{b}_2 \\ \mathbf{0}^\top & 1 \end{pmatrix} = \begin{pmatrix} A_1 A_2 & \mathbf{b}_1 + A_1 \mathbf{b}_2 \\ \mathbf{0}^\top & 1 \end{pmatrix} \text{ and } \begin{pmatrix} A & \mathbf{b} \\ \mathbf{0}^\top & 1 \end{pmatrix}^{-1} = \begin{pmatrix} A^{-1} & -A^{-1} \mathbf{b} \\ \mathbf{0}^\top & 1 \end{pmatrix}$$

A group representation is a group homomorphism $\pi:G\to \mathrm{GL}(V)$, (i.e. $\pi(gh)=\pi(g)\pi(h)$) so that any element g can be seen as a matrix $\pi(g)$.



- (V, +). Any vector space V with vector addition as the binary operation.
- $(\mathbb{R}_{>0},\times)$. Positive real numbers $\mathbb{R}_{>0}$ under multiplication.
- $\mathfrak{GL}(n,\mathbb{R})$. Invertible $n\times n$ matrices, under matrix multiplication.
- \bullet Aff $(n,\mathbb{R})\cong \mathrm{GL}(n,\mathbb{R})\ltimes\mathbb{R}^n$. The Affine group, consisting of pairs (A,\mathbf{b}) where the group multiplication is $(A_1, \mathbf{b}_1)(A_2, \mathbf{b}_2) = (A_1A_2, A\mathbf{b}_2 + \mathbf{b}_1)$. Another way to realize this group is as a subgroup of $(n+1) \times (n+1)$ matrices.

$$\begin{pmatrix} A_1 & \mathbf{b}_1 \\ \mathbf{0}^\top & 1 \end{pmatrix} \begin{pmatrix} A_2 & \mathbf{b}_2 \\ \mathbf{0}^\top & 1 \end{pmatrix} = \begin{pmatrix} A_1 A_2 & \mathbf{b}_1 + A_1 \mathbf{b}_2 \\ \mathbf{0}^\top & 1 \end{pmatrix} \text{ and } \begin{pmatrix} A & \mathbf{b} \\ \mathbf{0}^\top & 1 \end{pmatrix}^{-1} = \begin{pmatrix} A^{-1} & -A^{-1} \mathbf{b} \\ \mathbf{0}^\top & 1 \end{pmatrix}$$

A group representation is a group homomorphism $\pi:G\to \mathrm{GL}(V)$, (i.e. $\pi(gh)=\pi(q)\pi(h)$) so that any element q can be seen as a matrix $\pi(q)$.

Elements of the symmetric group S_3 can be represented as matrices: $\pi((12)(3)) = \begin{pmatrix} 0 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \end{pmatrix}$.

Group Actions

A group G acts on a set X if there is a map

$$\cdot:G\times X\to X$$

satisfying

- $e \cdot x = x$ for all $x \in X$ where $e \in G$ is the identity.

Group Actions

A group G acts on a set X if there is a map

$$\cdot:G\times X\to X$$

satisfying

- $e \cdot x = x$ for all $x \in X$ where $e \in G$ is the identity.

Every group acts on itself with the group multiplication. In other words X=G and $g\cdot x=gx$.

Group Actions

A group G acts on a set X if there is a map

$$\cdot: G \times X \to X$$

satisfying

- $e \cdot x = x$ for all $x \in X$ where $e \in G$ is the identity.

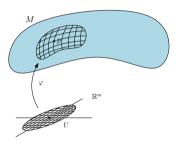
Every group acts on itself with the group multiplication. In other words X=G and $g\cdot x=gx$.

Another example, $(A, \mathbf{b}) \in \mathrm{Aff}(n, \mathbb{R})$ acting on $\mathbf{x} \in \mathbb{R}^n$ by

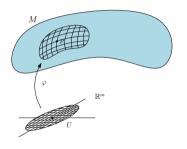
$$(A, \mathbf{b}) \cdot \mathbf{x} = A\mathbf{x} + \mathbf{b}.$$

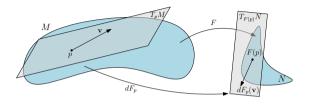


Manifolds are geometric objects which have local coordinates via charts $\varphi:U\subseteq\mathbb{R}^m\to M$.



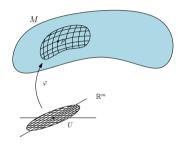
Manifolds are geometric objects which have local coordinates via charts $\varphi:U\subseteq\mathbb{R}^m\to M$.

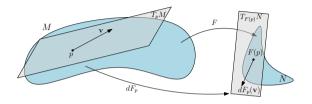




Differential of maps between two manifolds are linear maps between their tangent spaces.

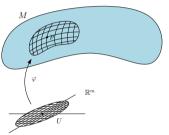
Manifolds are geometric objects which have local coordinates via charts $\varphi:U\subseteq\mathbb{R}^m\to M$.

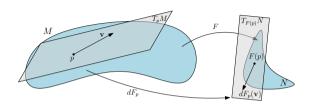




Differential of maps between two manifolds are linear maps between their tangent spaces. So if $N\subseteq\mathbb{R}$ —as $T_{f(p)}N\cong\mathbb{R}$ for any point p—we get a linear map $df_p:T_pM\to\mathbb{R}$, i.e. a linear functional: $df_p\in T_p^*M$.

Manifolds are geometric objects which have local coordinates via charts $\varphi:U\subseteq\mathbb{R}^m\to M$.





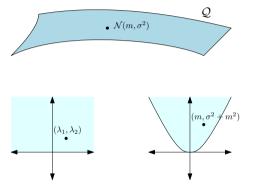
Differential of maps between two manifolds are linear maps between their tangent spaces. So if $N\subseteq\mathbb{R}$ —as $T_{f(p)}N\cong\mathbb{R}$ for any point p—we get a linear map $df_p:T_pM\to\mathbb{R}$, i.e. a linear functional: $df_p\in T_p^*M$.

If we have two maps $F:M\to N$ and $G:N\to P$ then the chain rule is

$$d(G \circ F)_p = dG_{F(p)}dF_p$$

On the right hand side, we have a composition of two linear maps, written multiplicatively.

Manifolds of distributions: more than one chart is possible

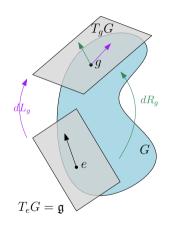


The 1-d Gaussian distributions form a two dimensional manifold with a single chart. But there can be more than one parametrization. One with natural parameters

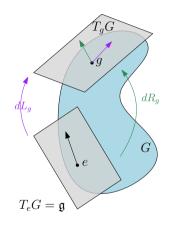
$$(\lambda_1, \lambda_2) \longmapsto q(x) \propto e^{-\lambda_1 x - \lambda_2 x^2}$$

or on the other hand, the expectations of x and x^2 are also enough to characterize $q(x) = \frac{1}{\sqrt{2\pi}\sigma}e^{-\frac{(x-m)^2}{2\sigma^2}}.$

The Euclidean inner product w.r.t. (λ_1, λ_2) or w.r.t. $(\mu_1, \mu_2) = (m, \sigma^2 + m^2)$ gives different angles between tangent vectors to $T_a Q$.

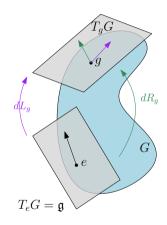


Lie groups are groups which are also manifolds. Multiplication and inversion maps are smooth.



Lie groups are groups which are also manifolds. Multiplication and inversion maps are smooth.

For every $g \in G$ left and right multiplication maps defined by $L_g(x) = gx, R_g(x) = xg$ send e to g.

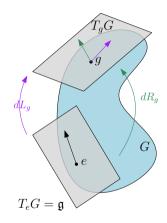


Lie groups are groups which are also manifolds.

Multiplication and inversion maps are smooth.

For every $g \in G$ left and right multiplication maps defined by $L_q(x) = gx, R_q(x) = xg$ send e to g.

Their differentials connect the tangent spaces \mathfrak{g} and T_aG .



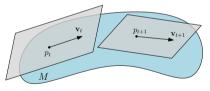
Lie groups are groups which are also manifolds.

Multiplication and inversion maps are smooth.

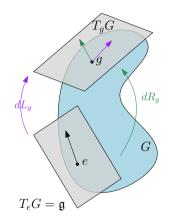
For every $g \in G$ left and right multiplication maps defined by $L_g(x) = gx, R_g(x) = xg$ send e to g.

Their differentials connect the tangent spaces $\mathfrak g$ and T_gG .

If learning on a general manifold, then vectors at different points are in totally different vector spaces. How to add them, like in momentum accumulation?



Lie groups: Multiplication connects tangent planes.



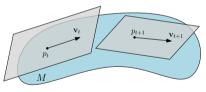
Lie groups are groups which are also manifolds.

Multiplication and inversion maps are smooth.

For every $g \in G$ left and right multiplication maps defined by $L_g(x) = gx, R_g(x) = xg$ send e to g.

Their differentials connect the tangent spaces $\mathfrak g$ and T_gG .

If learning on a general manifold, then vectors at different points are in totally different vector spaces. How to add them, like in momentum accumulation?

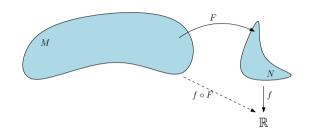


 $R_{g^{-1}} \circ L_g : e \mapsto e$, its differential is the adjoint map $\mathrm{Ad}_g : \mathfrak{g} \to \mathfrak{g}$. Identity for commutative groups.

Pushforwards and Pullbacks

Given $F:M\to N$, a map between manifolds, for every (real valued) function $f:N\to\mathbb{R}$ we can pull it back to a function $f\circ F:M\to\mathbb{R}$.

$$F^*: \mathcal{C}^{\infty}(N, \mathbb{R}) \longrightarrow \mathcal{C}^{\infty}(M, \mathbb{R})$$
$$f \longmapsto F^*(f) = f \circ F.$$

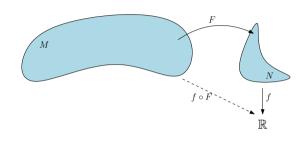


Pushforwards and Pullbacks

Given $F:M\to N$, a map between manifolds, for every (real valued) function $f:N\to\mathbb{R}$ we can pull it back to a function $f\circ F:M\to\mathbb{R}$.

$$F^*: \mathcal{C}^{\infty}(N, \mathbb{R}) \longrightarrow \mathcal{C}^{\infty}(M, \mathbb{R})$$

 $f \longmapsto F^*(f) = f \circ F.$



On the flip side probability distributions (or measures) μ on M are pushed to $F_*\mu$ on N:

$$\int_N f d(F_* \mu) = \int_M F^*(f) d\mu.$$

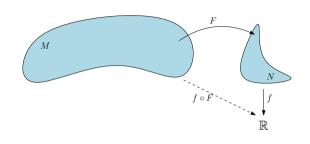


Pushforwards and Pullbacks

Given $F:M\to N$, a map between manifolds, for every (real valued) function $f:N\to\mathbb{R}$ we can pull it back to a function $f\circ F:M\to\mathbb{R}$.

$$F^*: \mathcal{C}^{\infty}(N, \mathbb{R}) \longrightarrow \mathcal{C}^{\infty}(M, \mathbb{R})$$

 $f \longmapsto F^*(f) = f \circ F.$



On the flip side probability distributions (or measures) μ on M are pushed to $F_*\mu$ on N:

$$\int_{N} f d(F_* \mu) = \int_{M} F^*(f) d\mu.$$

Another way to see this is on a measurable set $E \subseteq N$ the measure of the pushforward distribution is $(F^*\mu)(E) = \mu(F^{-1}(E)) = \mu(\{p \in M : F(p) \in E\}).$

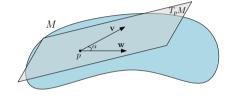


Pullbacks and pushforwards, II

The tangent vectors are pushed out to the range by the differential, a linear map

$$\mathrm{d}F_p:T_pM\to T_{F(p)}N.$$

Metrics, are bilinear forms on tangent spaces $\omega:T_pM\otimes T_pM\to\mathbb{R}$ which measure the length of the vectors and the angle between the vectors, just like the Euclidean dot product.

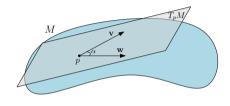


Pullbacks and pushforwards, II

The tangent vectors are pushed out to the range by the differential, a linear map

$$dF_p: T_pM \to T_{F(p)}N.$$

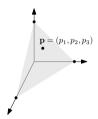
Metrics, are bilinear forms on tangent spaces $\omega:T_pM\otimes T_pM\to\mathbb{R}$ which measure the length of the vectors and the angle between the vectors, just like the Euclidean dot product.



A Riemannian manifold is a choice of bilinear form at every tangent plane T_pM on M. The vectors go forward, and a metric ω on N can be pulled back via

$$F^*(\omega)(\mathbf{v}, \mathbf{w}) = \omega(\mathrm{d}F_p\mathbf{v}, \mathrm{d}F_p\mathbf{w}).$$

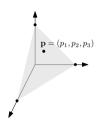




The group $(\mathbb{R}^r,+)$ acts on finite probability distributions on r points by

$$\mathbf{x} \cdot \mathbf{p} = \text{Norm}(e^{\mathbf{x}} \odot \mathbf{p}) = \left(\frac{e^{x_i} p_i}{\sum_j e^{x_j} p_j}\right)_{i=1,\dots,r}$$

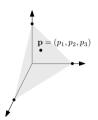




Assume G acts transitively on $\theta \in \Theta$. For example on the space of weights, e.g. \mathbb{R}^P .

The group $(\mathbb{R}^r,+)$ acts on finite probability distributions on r points by

$$\mathbf{x} \cdot \mathbf{p} = \text{Norm}(e^{\mathbf{x}} \odot \mathbf{p}) = \left(\frac{e^{x_i} p_i}{\sum_j e^{x_j} p_j}\right)_{i=1,\dots,r}$$

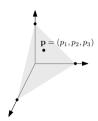


The group $(\mathbb{R}^r, +)$ acts on finite probability distributions on r points by

$$\mathbf{x} \cdot \mathbf{p} = \text{Norm}(e^{\mathbf{x}} \odot \mathbf{p}) = \left(\frac{e^{x_i} p_i}{\sum_j e^{x_j} p_j}\right)_{i=1,\dots,r}$$

Assume G acts transitively on $\theta \in \Theta$. For example on the space of weights, e.g. \mathbb{R}^P .

Every $a \in G$ defines a map $\theta \mapsto a \cdot \theta$.



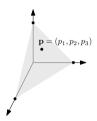
The group $(\mathbb{R}^r,+)$ acts on finite probability distributions on r points by

$$\mathbf{x} \cdot \mathbf{p} = \text{Norm}(e^{\mathbf{x}} \odot \mathbf{p}) = \left(\frac{e^{x_i} p_i}{\sum_j e^{x_j} p_j}\right)_{i=1,\dots,r}$$

Assume G acts transitively on $\theta \in \Theta$. For example on the space of weights, e.g. \mathbb{R}^P .

Every $g \in G$ defines a map $\theta \mapsto g \cdot \theta$.

$$u$$
 a base measure, s.t. $u(g\cdot E)=\chi(g)\nu(E)$, e.g. $\chi(g)=\det(A)$ for $g=(A,\mathbf{b})\in \mathrm{Aff}(n,\mathbb{R})$.



The group $(\mathbb{R}^r, +)$ acts on finite probability distributions on r points by

$$\mathbf{x} \cdot \mathbf{p} = \text{Norm}(e^{\mathbf{x}} \odot \mathbf{p}) = \left(\frac{e^{x_i} p_i}{\sum_j e^{x_j} p_j}\right)_{i=1,\dots,r}$$

Assume G acts transitively on $\theta \in \Theta$. For example on the space of weights, e.g. \mathbb{R}^P .

Every $g \in G$ defines a map $\theta \mapsto g \cdot \theta$.

$$u$$
 a base measure, s.t. $u(g\cdot E)=\chi(g)\nu(E),$ e.g. $\chi(g)=\det(A)$ for $g=(A,\mathbf{b})\in \mathrm{Aff}(n,\mathbb{R}).$

The pushforward of distributions $q\mathrm{d}\nu$ are given by $g_*(q\mathrm{d}\nu)=q^g\mathrm{d}\nu$ where

$$q^g(\boldsymbol{\theta}) = \frac{1}{\chi(g)} q(g^{-1} \cdot \boldsymbol{\theta}).$$



Assume a Lie group ${\cal G}$ acts on the parameter manifold Θ ,

Assume a Lie group G acts on the parameter manifold Θ , it also acts on distributions on Θ .

Assume a Lie group G acts on the parameter manifold Θ , it also acts on distributions on Θ . Q is formed as the orbit of such an action for any base distribution q_0 :

$$\mathcal{Q} = \{q^g \mathrm{d}\nu : g \in G\}. \qquad \text{where recall} \qquad q^g(\boldsymbol{\theta}) = \frac{1}{\chi(q)} q_0(g^{-1} \cdot \boldsymbol{\theta}).$$

Assume a Lie group G acts on the parameter manifold Θ , it also acts on distributions on Θ . Q is formed as the orbit of such an action for any base distribution q_0 :

$$Q = \{q^g \operatorname{d}\nu : g \in G\}. \qquad \text{when}$$

$$q_0 = \operatorname{translation}_{g} \qquad q^g(\theta) = q_0(\theta - g)$$

$$q_0 = \operatorname{translation}_{g} \qquad q_0(\theta) = \operatorname{$$

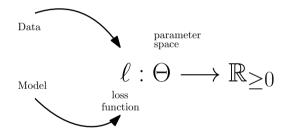
where recall
$$q^g(\boldsymbol{\theta}) = \frac{1}{\chi(g)} q_0(g^{-1} \cdot \boldsymbol{\theta}).$$

$$ullet$$
 $G=(\mathbb{R},+)$, $\Theta=\mathbb{R}$,

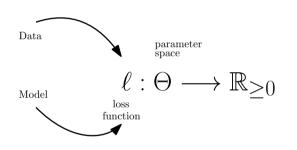
•
$$G = (\mathbb{R}_{>0}, \times), \Theta = \mathbb{R}_{>0},$$

•
$$G = \mathrm{Aff}(\mathbb{R}) = \mathbb{R}_{>0} \ltimes \mathbb{R}, \ \Theta = \mathbb{R}$$

The classical and Bayesian learning setups (what we learn!)



The classical and Bayesian learning setups (what we learn!)



Classically: find $\theta^* \in \Theta$ minimizing ℓ .

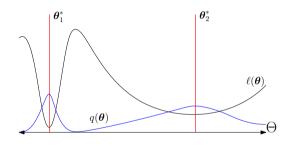
Bayesian : find a distribution $q \in \mathcal{P}_{\nu}(\Theta)$

Classical vs. Bayesian learning

The loss function is highly nonconvex. Usually

$$\ell(\boldsymbol{\theta}) = \sum_{i=1}^{N} \ell_i(\boldsymbol{\theta}) + R(\boldsymbol{\theta})$$

where $\ell_i(\boldsymbol{\theta})$ is the loss contribution from the i^{th} data point and $R(\boldsymbol{\theta})$ regularizer.



 θ_1^* and θ_2^* are both equally valid explanations of the same data. A distribution over the data considers both explanations "at the same time".

Betting it all on one outcome

Say two dice are thrown and I tell you that the sum is greater than 7.

Betting it all on one outcome

Say two dice are thrown and I tell you that the sum is greater than $7. \odot, \odot$ satisfies this.

We could say the result was definitely . . .

Betting it all on one outcome

Say two dice are thrown and I tell you that the sum is greater than 7.[™], [™] satisfies this.

We could say the result was definitely **≅**, **⊙**.

But there are a total of 15 possibilities

It is much more sensible to say it is one of these 15 outcomes, with equal probability. (principle of indifference, principle of maximum entropy)

 $\ell(\theta)$, a loss function on model parameters $\theta \in \Theta$. Pick a base measure ν on Θ .

 $\ell(\theta)$, a loss function on model parameters $\theta \in \Theta$. Pick a base measure ν on Θ . We solve

$$q_* \in \underset{q \in \mathcal{Q}}{\operatorname{arg min}} \underbrace{\mathbb{E}_q[\ell] - \tau \mathcal{H}_{\nu}(q)}_{=:\mathcal{E}(q)}$$

for some family of distributions $Q \subseteq \mathcal{P}_{\nu}(\Theta) = \{q(\boldsymbol{\theta}) d\nu(\boldsymbol{\theta})\}$ on the parameters.

 $\ell(\theta)$, a loss function on model parameters $\theta \in \Theta$. Pick a base measure ν on Θ . We solve

$$q_* \in \underset{q \in \mathcal{Q}}{\operatorname{arg min}} \underbrace{\mathbb{E}_q[\ell] - \tau \mathcal{H}_\nu(q)}_{=:\mathcal{E}(q)}$$

for some family of distributions $Q \subseteq \mathcal{P}_{\nu}(\Theta) = \{q(\theta) d\nu(\theta)\}$ on the parameters.

• The expectation $\mathbb{E}_q[\ell] = \int_{\Theta} \ell(\boldsymbol{\theta}) q(\boldsymbol{\theta}) d\nu(\boldsymbol{\theta})$ prefers regions with low loss.

 $\ell(\boldsymbol{\theta})$, a loss function on model parameters $\boldsymbol{\theta} \in \Theta$. Pick a base measure ν on Θ . We solve

$$q_* \in \underset{q \in \mathcal{Q}}{\operatorname{arg min}} \underbrace{\mathbb{E}_q[\ell] - \tau \mathcal{H}_{\nu}(q)}_{=:\mathcal{E}(q)}$$

for some family of distributions $Q \subseteq \mathcal{P}_{\nu}(\Theta) = \{q(\theta) d\nu(\theta)\}$ on the parameters.

- The expectation $\mathbb{E}_q[\ell] = \int_{\Theta} \ell(\boldsymbol{\theta}) q(\boldsymbol{\theta}) \mathrm{d}\nu(\boldsymbol{\theta})$ prefers regions with low loss.
- The entropy $\mathcal{H}_{\nu}(q) = -\int_{\Theta} q(\boldsymbol{\theta}) \log q(\boldsymbol{\theta}) \mathrm{d}\nu(\boldsymbol{\theta})$ prefers a higher spread of q.

 $\ell(\boldsymbol{\theta})$, a loss function on model parameters $\boldsymbol{\theta} \in \Theta$. Pick a base measure ν on Θ . We solve

$$q_* \in \underset{q \in \mathcal{Q}}{\operatorname{arg min}} \underbrace{\mathbb{E}_q[\ell] - \tau \mathcal{H}_{\nu}(q)}_{=:\mathcal{E}(q)}$$

for some family of distributions $Q \subseteq \mathcal{P}_{\nu}(\Theta) = \{q(\theta) d\nu(\theta)\}$ on the parameters.

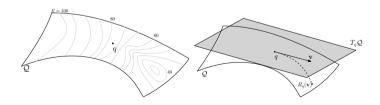
- The expectation $\mathbb{E}_q[\ell] = \int_{\Theta} \ell(\boldsymbol{\theta}) q(\boldsymbol{\theta}) \mathrm{d}\nu(\boldsymbol{\theta})$ prefers regions with low loss.
- The entropy $\mathcal{H}_{\nu}(q) = -\int_{\Theta} q(\boldsymbol{\theta}) \log q(\boldsymbol{\theta}) d\nu(\boldsymbol{\theta})$ prefers a higher spread of q.
- The temperature $\tau > 0$ is a balancing term.



Gradient Descent on Manifolds

Gradient descent for $\mathcal{E}(q)$ on the manifold $\mathcal Q$ has three components.

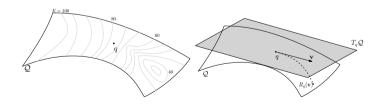
- **①** Calculate the differential $d\mathcal{E}|_q: T_q\mathcal{Q} \mapsto \mathbb{R}$.
- ② Find a way to turn the covector $d\mathcal{E}|_q \in T_q^*\mathcal{Q}$ into a vector $\mathbf{v} \in T_q\mathcal{Q}$.
- **3** A choice of retraction brings us back onto the manifold $R_q(\mathbf{v}) \in \mathcal{Q}$.

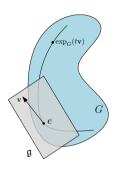


Gradient Descent on Manifolds

Gradient descent for $\mathcal{E}(q)$ on the manifold $\mathcal Q$ has three components.

- Calculate the differential $d\mathcal{E}|_q:T_q\mathcal{Q}\mapsto\mathbb{R}.$
- ② Find a way to turn the covector $d\mathcal{E}|_q \in T_q^*\mathcal{Q}$ into a vector $\mathbf{v} \in T_q\mathcal{Q}$.
- lack A choice of retraction brings us back onto the manifold $R_q(\mathbf{v}) \in \mathcal{Q}.$





On every Lie group there is an exponential map

$$\exp_G: \mathfrak{g} \to G$$

provides a canonical retraction.

We now solve

$$\underset{g \in G}{\arg\min} \, \mathcal{E}(q^g) = \underset{g \in G}{\arg\min} \int_{\Theta} q^g(\boldsymbol{\theta}) \ell(\boldsymbol{\theta}) + \tau q^g(\boldsymbol{\theta}) \log q^g(\boldsymbol{\theta}) \mathrm{d}\nu(\boldsymbol{\theta})$$

Given $X \in \mathfrak{g} = T_eG$ the differential in the direction of X is

$$\frac{\mathrm{d}}{\mathrm{d}t}\mathcal{E}(q^{ge^{tX}})\bigg|_{t=0} = \frac{\mathrm{d}}{\mathrm{d}t}\underbrace{\int_{\Theta}q^{ge^{tX}}(\boldsymbol{\theta})\ell(\boldsymbol{\theta})\mathrm{d}\nu(\boldsymbol{\theta})}_{\text{data contribution}} + \underbrace{\tau\int_{\Theta}q^{ge^{tX}}(\boldsymbol{\theta})\log q^{ge^{tX}}(\boldsymbol{\theta})\mathrm{d}\nu(\boldsymbol{\theta})}_{\text{entropy contribution}}\bigg|_{t=0}$$

We now solve

$$\underset{g \in G}{\arg\min} \, \mathcal{E}(q^g) = \underset{g \in G}{\arg\min} \int_{\Theta} q^g(\boldsymbol{\theta}) \ell(\boldsymbol{\theta}) + \tau q^g(\boldsymbol{\theta}) \log q^g(\boldsymbol{\theta}) \mathrm{d}\nu(\boldsymbol{\theta})$$

Given $X \in \mathfrak{g} = T_eG$ the differential in the direction of X is

$$\frac{\mathrm{d}}{\mathrm{d}t}\mathcal{E}(q^{ge^{tX}})\bigg|_{t=0} = \frac{\mathrm{d}}{\mathrm{d}t}\underbrace{\int_{\Theta}q^{ge^{tX}}(\boldsymbol{\theta})\ell(\boldsymbol{\theta})\mathrm{d}\nu(\boldsymbol{\theta})}_{\text{data contribution}} + \underbrace{\tau\int_{\Theta}q^{ge^{tX}}(\boldsymbol{\theta})\log q^{ge^{tX}}(\boldsymbol{\theta})\mathrm{d}\nu(\boldsymbol{\theta})}_{\text{entropy contribution}}\bigg|_{t=0}$$

$$\left. \frac{\mathrm{d}}{\mathrm{d}t} \int_{\Theta} q^{ge^{tX}}(\boldsymbol{\theta}) \ell(\boldsymbol{\theta}) \mathrm{d}\nu(\boldsymbol{\theta}) \right|_{t=0}$$



We now solve

$$\underset{g \in G}{\arg\min} \, \mathcal{E}(q^g) = \underset{g \in G}{\arg\min} \int_{\Theta} q^g(\boldsymbol{\theta}) \ell(\boldsymbol{\theta}) + \tau q^g(\boldsymbol{\theta}) \log q^g(\boldsymbol{\theta}) \mathrm{d}\nu(\boldsymbol{\theta})$$

Given $X \in \mathfrak{g} = T_eG$ the differential in the direction of X is

$$\frac{\mathrm{d}}{\mathrm{d}t}\mathcal{E}(q^{ge^{tX}})\bigg|_{t=0} = \frac{\mathrm{d}}{\mathrm{d}t}\underbrace{\int_{\Theta}q^{ge^{tX}}(\boldsymbol{\theta})\ell(\boldsymbol{\theta})\mathrm{d}\nu(\boldsymbol{\theta})}_{\text{data contribution}} + \underbrace{\tau\int_{\Theta}q^{ge^{tX}}(\boldsymbol{\theta})\log q^{ge^{tX}}(\boldsymbol{\theta})\mathrm{d}\nu(\boldsymbol{\theta})}_{\text{entropy contribution}}\bigg|_{t=0}$$

$$\frac{\mathrm{d}}{\mathrm{d}t} \int_{\Theta} \frac{1}{\chi(ge^{tX})} q_0((ge^{tX})^{-1} \cdot \boldsymbol{\theta}) \ell(\boldsymbol{\theta}) \mathrm{d}\nu(\boldsymbol{\theta}) \bigg|_{t=0}$$



We now solve

$$\underset{g \in G}{\arg\min} \, \mathcal{E}(q^g) = \underset{g \in G}{\arg\min} \int_{\Theta} q^g(\boldsymbol{\theta}) \ell(\boldsymbol{\theta}) + \tau q^g(\boldsymbol{\theta}) \log q^g(\boldsymbol{\theta}) \mathrm{d}\nu(\boldsymbol{\theta})$$

Given $X \in \mathfrak{g} = T_eG$ the differential in the direction of X is

$$\frac{\mathrm{d}}{\mathrm{d}t}\mathcal{E}(q^{ge^{tX}})\bigg|_{t=0} = \frac{\mathrm{d}}{\mathrm{d}t}\underbrace{\int_{\Theta}q^{ge^{tX}}(\boldsymbol{\theta})\ell(\boldsymbol{\theta})\mathrm{d}\nu(\boldsymbol{\theta})}_{\text{data contribution}} + \underbrace{\tau\int_{\Theta}q^{ge^{tX}}(\boldsymbol{\theta})\log q^{ge^{tX}}(\boldsymbol{\theta})\mathrm{d}\nu(\boldsymbol{\theta})}_{\text{entropy contribution}}\bigg|_{t=0}$$

$$\frac{\mathrm{d}}{\mathrm{d}t} \int_{\Theta} \frac{1}{\chi(ge^{tX})} q_0(\boldsymbol{\theta}) \ell(ge^{tX} \cdot \boldsymbol{\theta}) \mathrm{d}\nu(ge^{tX} \cdot \boldsymbol{\theta}) \bigg|_{t=0}$$



We now solve

$$\underset{g \in G}{\arg\min} \, \mathcal{E}(q^g) = \underset{g \in G}{\arg\min} \int_{\Theta} q^g(\boldsymbol{\theta}) \ell(\boldsymbol{\theta}) + \tau q^g(\boldsymbol{\theta}) \log q^g(\boldsymbol{\theta}) \mathrm{d}\nu(\boldsymbol{\theta})$$

Given $X \in \mathfrak{g} = T_eG$ the differential in the direction of X is

$$\frac{\mathrm{d}}{\mathrm{d}t}\mathcal{E}(q^{ge^{tX}})\bigg|_{t=0} = \frac{\mathrm{d}}{\mathrm{d}t}\underbrace{\int_{\Theta}q^{ge^{tX}}(\boldsymbol{\theta})\ell(\boldsymbol{\theta})\mathrm{d}\nu(\boldsymbol{\theta})}_{\text{data contribution}} + \underbrace{\tau\int_{\Theta}q^{ge^{tX}}(\boldsymbol{\theta})\log q^{ge^{tX}}(\boldsymbol{\theta})\mathrm{d}\nu(\boldsymbol{\theta})}_{\text{entropy contribution}}\bigg|_{t=0}$$

$$\frac{\mathrm{d}}{\mathrm{d}t} \int_{\Theta} q_0(\boldsymbol{\theta}) \ell(g e^{tX} \cdot \boldsymbol{\theta}) \mathrm{d}\nu(\boldsymbol{\theta}) \bigg|_{t=0}$$



We now solve

$$\underset{g \in G}{\arg\min} \, \mathcal{E}(q^g) = \underset{g \in G}{\arg\min} \int_{\Theta} q^g(\boldsymbol{\theta}) \ell(\boldsymbol{\theta}) + \tau q^g(\boldsymbol{\theta}) \log q^g(\boldsymbol{\theta}) \mathrm{d}\nu(\boldsymbol{\theta})$$

Given $X \in \mathfrak{g} = T_eG$ the differential in the direction of X is

$$\frac{\mathrm{d}}{\mathrm{d}t}\mathcal{E}(q^{ge^{tX}})\bigg|_{t=0} = \frac{\mathrm{d}}{\mathrm{d}t}\underbrace{\int_{\Theta}q^{ge^{tX}}(\boldsymbol{\theta})\ell(\boldsymbol{\theta})\mathrm{d}\nu(\boldsymbol{\theta})}_{\text{data contribution}} + \underbrace{\tau\int_{\Theta}q^{ge^{tX}}(\boldsymbol{\theta})\log q^{ge^{tX}}(\boldsymbol{\theta})\mathrm{d}\nu(\boldsymbol{\theta})}_{\text{entropy contribution}}\bigg|_{t=0}$$

$$\frac{\mathrm{d}}{\mathrm{d}t} \int_{\Theta} q_0(\boldsymbol{\theta}) \ell(g e^{tX} g^{-1} g \cdot \boldsymbol{\theta}) \mathrm{d}\nu(\boldsymbol{\theta}) \bigg|_{t=0}$$



We now solve

$$\underset{g \in G}{\arg\min} \, \mathcal{E}(q^g) = \underset{g \in G}{\arg\min} \int_{\Theta} q^g(\boldsymbol{\theta}) \ell(\boldsymbol{\theta}) + \tau q^g(\boldsymbol{\theta}) \log q^g(\boldsymbol{\theta}) \mathrm{d}\nu(\boldsymbol{\theta})$$

Given $X \in \mathfrak{g} = T_eG$ the differential in the direction of X is

$$\frac{\mathrm{d}}{\mathrm{d}t}\mathcal{E}(q^{ge^{tX}})\bigg|_{t=0} = \frac{\mathrm{d}}{\mathrm{d}t}\underbrace{\int_{\Theta}q^{ge^{tX}}(\boldsymbol{\theta})\ell(\boldsymbol{\theta})\mathrm{d}\nu(\boldsymbol{\theta})}_{\text{data contribution}} + \underbrace{\tau\int_{\Theta}q^{ge^{tX}}(\boldsymbol{\theta})\log q^{ge^{tX}}(\boldsymbol{\theta})\mathrm{d}\nu(\boldsymbol{\theta})}_{\text{entropy contribution}}\bigg|_{t=0}$$

$$\frac{\mathrm{d}}{\mathrm{d}t} \int_{\Theta} \frac{1}{\chi(g)} q_0(g^{-1} \cdot \boldsymbol{\theta}) \ell(g e^{tX} g^{-1} \cdot \boldsymbol{\theta}) \mathrm{d}\nu(\boldsymbol{\theta}) \bigg|_{t=0}$$



We now solve

$$\underset{g \in G}{\arg\min} \, \mathcal{E}(q^g) = \underset{g \in G}{\arg\min} \int_{\Theta} q^g(\boldsymbol{\theta}) \ell(\boldsymbol{\theta}) + \tau q^g(\boldsymbol{\theta}) \log q^g(\boldsymbol{\theta}) \mathrm{d}\nu(\boldsymbol{\theta})$$

Given $X \in \mathfrak{g} = T_eG$ the differential in the direction of X is

$$\frac{\mathrm{d}}{\mathrm{d}t}\mathcal{E}(q^{ge^{tX}})\bigg|_{t=0} = \frac{\mathrm{d}}{\mathrm{d}t}\underbrace{\int_{\Theta}q^{ge^{tX}}(\boldsymbol{\theta})\ell(\boldsymbol{\theta})\mathrm{d}\nu(\boldsymbol{\theta})}_{\text{data contribution}} + \underbrace{\tau\int_{\Theta}q^{ge^{tX}}(\boldsymbol{\theta})\log q^{ge^{tX}}(\boldsymbol{\theta})\mathrm{d}\nu(\boldsymbol{\theta})}_{\text{entropy contribution}}\bigg|_{t=0}$$

$$\frac{\mathrm{d}}{\mathrm{d}t} \int_{\Theta} q^{g}(\boldsymbol{\theta}) \ell(g e^{tX} g^{-1} \cdot \boldsymbol{\theta}) \mathrm{d}\nu(\boldsymbol{\theta}) \bigg|_{t=0}$$



Optimization on the group: calculating the differential

We now solve

$$\underset{g \in G}{\arg\min} \, \mathcal{E}(q^g) = \underset{g \in G}{\arg\min} \int_{\Theta} q^g(\boldsymbol{\theta}) \ell(\boldsymbol{\theta}) + \tau q^g(\boldsymbol{\theta}) \log q^g(\boldsymbol{\theta}) \mathrm{d}\nu(\boldsymbol{\theta})$$

Given $X \in \mathfrak{g} = T_eG$ the differential in the direction of X is

$$\frac{\mathrm{d}}{\mathrm{d}t}\mathcal{E}(q^{ge^{tX}})\bigg|_{t=0} = \frac{\mathrm{d}}{\mathrm{d}t}\underbrace{\int_{\Theta}q^{ge^{tX}}(\boldsymbol{\theta})\ell(\boldsymbol{\theta})\mathrm{d}\nu(\boldsymbol{\theta})}_{\text{data contribution}} + \underbrace{\tau\int_{\Theta}q^{ge^{tX}}(\boldsymbol{\theta})\log q^{ge^{tX}}(\boldsymbol{\theta})\mathrm{d}\nu(\boldsymbol{\theta})}_{\text{entropy contribution}}\bigg|_{t=0}$$

The data contribution can be rewritten as

$$\int_{\Theta} q^{g}(\boldsymbol{\theta}) \nabla \ell (g e^{\mathbf{0}} g^{-1} \cdot \boldsymbol{\theta})^{\top} \frac{\mathrm{d}}{\mathrm{d}t} \left((g e^{tX} g^{-1}) \cdot \boldsymbol{\theta} \right) \Big|_{t=0} \mathrm{d}\nu(\boldsymbol{\theta})$$



Optimization on the group: calculating the differential

We now solve

$$\underset{g \in G}{\arg\min} \, \mathcal{E}(q^g) = \underset{g \in G}{\arg\min} \int_{\Theta} q^g(\boldsymbol{\theta}) \ell(\boldsymbol{\theta}) + \tau q^g(\boldsymbol{\theta}) \log q^g(\boldsymbol{\theta}) \mathrm{d}\nu(\boldsymbol{\theta})$$

Given $X \in \mathfrak{g} = T_eG$ the differential in the direction of X is

$$\frac{\mathrm{d}}{\mathrm{d}t}\mathcal{E}(q^{ge^{tX}})\bigg|_{t=0} = \frac{\mathrm{d}}{\mathrm{d}t}\underbrace{\int_{\Theta}q^{ge^{tX}}(\boldsymbol{\theta})\ell(\boldsymbol{\theta})\mathrm{d}\nu(\boldsymbol{\theta})}_{\text{data contribution}} + \underbrace{\tau\int_{\Theta}q^{ge^{tX}}(\boldsymbol{\theta})\log q^{ge^{tX}}(\boldsymbol{\theta})\mathrm{d}\nu(\boldsymbol{\theta})}_{\text{entropy contribution}}\bigg|_{t=0}$$

The data contribution can be rewritten as

$$\int_{\Theta} q^g(\boldsymbol{\theta}) \nabla \ell(\boldsymbol{\theta})^{\top} (\mathrm{Ad}_g(X) \cdot \boldsymbol{\theta}) \mathrm{d}\nu(\boldsymbol{\theta})$$



Optimization on the group: calculating the differential

We now solve

$$\underset{g \in G}{\arg\min} \, \mathcal{E}(q^g) = \underset{g \in G}{\arg\min} \int_{\Theta} q^g(\boldsymbol{\theta}) \ell(\boldsymbol{\theta}) + \tau q^g(\boldsymbol{\theta}) \log q^g(\boldsymbol{\theta}) \mathrm{d}\nu(\boldsymbol{\theta})$$

Given $X \in \mathfrak{g} = T_eG$ the differential in the direction of X is

$$\frac{\mathrm{d}}{\mathrm{d}t}\mathcal{E}(q^{ge^{tX}})\bigg|_{t=0} = \frac{\mathrm{d}}{\mathrm{d}t}\underbrace{\int_{\Theta}q^{ge^{tX}}(\boldsymbol{\theta})\ell(\boldsymbol{\theta})\mathrm{d}\nu(\boldsymbol{\theta})}_{\text{data contribution}} + \underbrace{\tau\int_{\Theta}q^{ge^{tX}}(\boldsymbol{\theta})\log q^{ge^{tX}}(\boldsymbol{\theta})\mathrm{d}\nu(\boldsymbol{\theta})}_{\text{entropy contribution}}\bigg|_{t=0}$$

The data contribution can be rewritten as

$$\int_{\Theta} q^{g}(\boldsymbol{\theta}) \nabla \ell(\boldsymbol{\theta})^{\top} (\mathrm{Ad}_{g}(X) \cdot \boldsymbol{\theta}) \mathrm{d}\nu(\boldsymbol{\theta}) \approx \frac{1}{K} \sum_{\substack{i=1 \\ \theta_{i} \sim q^{g}}}^{K} \nabla \ell(\boldsymbol{\theta}_{i})^{\top} (\mathrm{Ad}_{g}(X) \cdot \boldsymbol{\theta}_{i})$$

The metric determines the fastest direction of descent.

Which descent direction is *fastest* depends on how we measure distances in the tangent space.

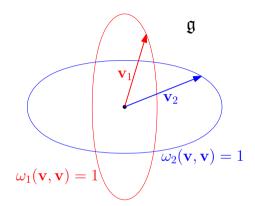
 \mathbf{v}_1 is fastest if we measure $\|v\|$ using ω_1

 \mathbf{v}_2 is fastest if $\|v\|$ is measured using $\omega_2.$

The pullback the Fisher information metric on $T_{q^g}\mathcal{Q}$

$$\omega^{\mathsf{Fisher}}(h_1, h_2) = \mathbb{E}_q \left[(\partial_{h_1} \log q^g) (\partial_{h_2} \log q^g) \right]$$

to \mathfrak{g} is independent of $g \in G$.



The metric determines the fastest direction of descent.

Which descent direction is *fastest* depends on how we measure distances in the tangent space.

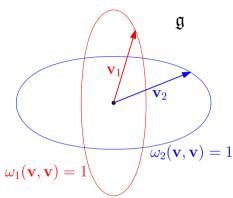
 \mathbf{v}_1 is fastest if we measure $\|v\|$ using ω_1

 \mathbf{v}_2 is fastest if $\|v\|$ is measured using $\omega_2.$

The pullback the Fisher information metric on $T_{q^g}\mathcal{Q}$

$$\omega^{\mathsf{Fisher}}(h_1, h_2) = \mathbb{E}_q \left[(\partial_{h_1} \log q^g) (\partial_{h_2} \log q^g) \right]$$

to \mathfrak{g} is independent of $g \in G$.

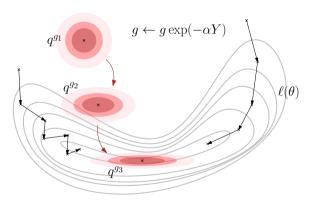


Indeed
$$h_X^g(\boldsymbol{\theta}) := \frac{\mathrm{d}}{\mathrm{d}t} q^{ge^{tX}}(\boldsymbol{\theta}) \Big|_{t=0} = \frac{1}{\chi(g)} h_X^e(g^{-1} \cdot \boldsymbol{\theta})$$
, so a change of variables $\boldsymbol{\theta} \mapsto g \cdot \boldsymbol{\theta}$ shows

$$\omega^{g}(X,Y) = \int_{\Theta} \frac{h_X^{g}(\boldsymbol{\theta})}{q^{g}(\boldsymbol{\theta})} \frac{h_X^{g}(\boldsymbol{\theta})}{q^{g}(\boldsymbol{\theta})} q^{g}(\boldsymbol{\theta}) d\nu(\boldsymbol{\theta}) = \omega^{e}(X,Y).$$

Classical Learning vs. Learning via Group

The *point based* gradient descent updates parameters: $\theta \leftarrow \theta - \alpha \nabla \ell(\theta)$ Bayesian Learning Rule(s) update the distribution over the parameters θ .



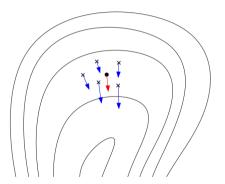
 $Y \in T_eG$ is the direction of fastest ascent of $\mathcal{E}(q^g)$ w.r.t. the Fisher metric.

Specific Update Formulas: The Additive Group

$$g \in \mathbb{R}^P$$
 additive \Longrightarrow $g \longleftarrow g - \alpha \mathbb{E}_{q_g} \big[\nabla_{\pmb{\theta}} \ell(\pmb{\theta}) \big]$

Specific Update Formulas: The Additive Group

$$g \in \mathbb{R}^P$$
 additive \Longrightarrow $g \longleftarrow g - \alpha \mathbb{E}_{q_g} \big[\nabla_{\pmb{\theta}} \ell(\pmb{\theta}) \big]$



Instead of going in the direction of the derivative at g, the direction is chosen by consensus with at points sampled from q_a .

Multiplicative and Affine Update Formulas

If $g \in (\mathbb{R}_{>0}^P, \times)$ acting on a parameters $\boldsymbol{\theta} \in \mathbb{R}_+^P$:

$$X = \mathbb{E}_{q^g}[\boldsymbol{\theta} \odot \nabla_{\boldsymbol{\theta}} \ell(\boldsymbol{\theta})] - \tau$$
$$g \longleftarrow g \odot e^{-\alpha X}$$

Multiplicative and Affine Update Formulas

If $g \in (\mathbb{R}^P_{>0}, \times)$ acting on a parameters $\boldsymbol{\theta} \in \mathbb{R}^P_+$:

$$X = \mathbb{E}_{q^g}[\boldsymbol{\theta} \odot \nabla_{\boldsymbol{\theta}} \ell(\boldsymbol{\theta})] - \tau$$
$$g \longleftarrow g \odot e^{-\alpha X}$$

For the Affine group the Lie algebra elements (tangent vectors) are of the form $\begin{pmatrix} X & \mathbf{y} \\ \mathbf{0}^\top & 0 \end{pmatrix}$ and the gradient update at q^g is

$$X = \mathbb{E}_{q_0} \left[A^{\top} \nabla_{\boldsymbol{\theta}} \ell (A \boldsymbol{\theta} + \mathbf{b}) \boldsymbol{\theta}^{\top} \right] - \tau I \qquad \mathbf{y} = \mathbb{E}_{q_0} \left[A^{\top} \nabla \ell (A \boldsymbol{\theta} + \mathbf{b}) \right].$$

with a slightly more involved exponential map

$$A \longleftarrow Ae^{-\alpha X}$$
 $\mathbf{b} \longleftarrow A\frac{e^{-\alpha X} - I}{X}\mathbf{y} + \mathbf{b}$



Filters of the multiplicative group

Label nodes in a neural network "excitatory" or "inhibitory" like biology.

Magnitudes of the weights (in $\mathbb{R}_{>0}$) are the parameters (signs are fixed).

At each layer the map is $\mathbf{x} \mapsto \sigma(W_+\mathbf{x} - W_-\mathbf{x})$.

Filters of the multiplicative group

Label nodes in a neural network "excitatory" or "inhibitory" like biology.

Magnitudes of the weights (in $\mathbb{R}_{>0}$) are the parameters (signs are fixed).

At each layer the map is $\mathbf{x} \mapsto \sigma(W_+\mathbf{x} - W_-\mathbf{x})$.

Given $g \in \mathbb{R}^P_{>0}$, and q_0 Rayleigh, say, and $\theta_j \sim q_0^P$ for $j=1,\ldots,K$

$$M \leftarrow \beta M + (1 - \beta) \frac{1}{K} \sum_{j=1}^{K} (g \odot \boldsymbol{\theta}_j) \nabla \ell(g \odot \boldsymbol{\theta}_j) - \tau$$
$$g \leftarrow g \odot \exp(-\alpha M)$$

Filters of the multiplicative group

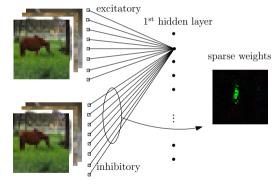
Label nodes in a neural network "excitatory" or "inhibitory" like biology.

Magnitudes of the weights (in $\mathbb{R}_{>0}$) are the parameters (signs are fixed).

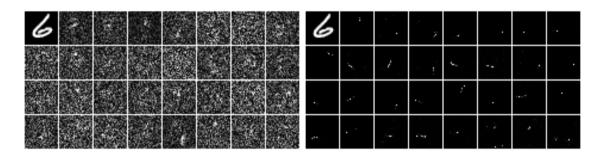
At each layer the map is $\mathbf{x} \mapsto \sigma(W_+\mathbf{x} - W_-\mathbf{x})$.

Given $g \in \mathbb{R}^P_{>0}$, and q_0 Rayleigh, say, and $\theta_j \sim q_0^P$ for $j=1,\dots,K$

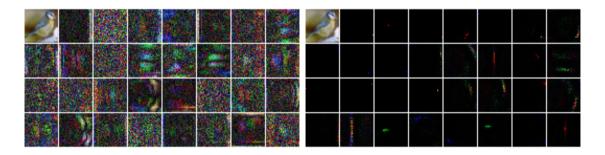
$$M \leftarrow \beta M + (1 - \beta) \frac{1}{K} \sum_{j=1}^{K} (g \odot \boldsymbol{\theta}_j) \nabla \ell(g \odot \boldsymbol{\theta}_j) - \tau$$
$$q \leftarrow q \odot \exp(-\alpha M)$$



Multiplicative vs Additive filters



The additive vs multiplicative filters for RGB images



Teşekkürler ありがとうございます Vielen Danke Merci Thank you.²

²More can be found at my blog-posts via https://ekiral.github.io/blog/blog-index.html

Stiefel Manifold Update

Assume parameters are given as a matrix and want to preserve orthogonality of columns.

$$\Theta = \operatorname{St}(n, m) = \{ \theta \in \operatorname{Mat}(n, m) : \theta^{\top} \theta = I_{m \times m} \}$$

The group $S = \mathrm{SO}(n)$ preserves this manifold. And given a loss function $\ell: \Theta \to \mathbb{R}_{\geq 0}$

$$Y \in \mathfrak{so}(n)$$
 the update direction $Y = \operatorname{Skew} Y_0 = \frac{Y_0 - Y_0^{ op}}{2}$ $Y_0 = \mathbb{E}_{q_{\Lambda}}[\nabla \ell \theta^{ op}]$

Stiefel Manifold Update

Assume parameters are given as a matrix and want to preserve orthogonality of columns.

$$\Theta = \operatorname{St}(n, m) = \{ \theta \in \operatorname{Mat}(n, m) : \theta^{\top} \theta = I_{m \times m} \}$$

The group S = SO(n) preserves this manifold. And given a loss function $\ell: \Theta \to \mathbb{R}_{\geq 0}$

$$Y \in \mathfrak{so}(n)$$
 the update direction $Y = \operatorname{Skew} Y_0 = \frac{Y_0 - Y_0^{ op}}{2}$ $Y_0 = \mathbb{E}_{q_{\Lambda}}[\nabla \ell \theta^{ op}]$

Here the distributions are parametrized by $\Lambda \in \mathrm{Mat}(n,m)$

$$q_{\Lambda}(\theta) \propto e^{-\operatorname{Tr}(\Lambda^{\top}\theta)}$$

and the update is given by

$$\Lambda \leftarrow e^{-\alpha Y} \Lambda$$
 (actually an efficient variation is used)



Koichi Tojo, Taro Yoshino's: "Harmonic Exponential Families".

G a Lie group $H \leq G$. Let ν be a relatively invariant measure on G $\pi: G \to \mathrm{GL}(V)$ a representation of G. Let α be a 1-cocycle of π such that $\alpha|_H \equiv 0$. So $\alpha: G \to V$ satisfies

$$\alpha(gh) = \pi(g)\alpha(h) + \alpha(g) = \alpha(g). \qquad \text{So } \alpha: \underbrace{G/H}_{:=\Theta} \to V$$

Let ν be a relatively invariant measure on Θ , meaning $\nu(gE)=\chi(g)\nu(E)$ for some homomorphism χ . Let $\lambda\in V^\vee$ s.t. $A(\lambda)=\log\int_\Theta e^{-\langle\lambda,\alpha(\theta)\rangle}\mathrm{d}\nu(\theta)<\infty$. For such λ

$$q_{\lambda}(\theta) d\nu(\theta) := e^{-\langle \lambda, \alpha(\theta) \rangle - A(\lambda)} d\nu(\theta)$$

forms an exponential family satisfying

$$\frac{1}{\chi(g)}q_{\lambda}(g^{-1}\theta) := q_{\pi^{\vee}(g)\lambda}(\theta) \qquad \text{where } \langle \pi^{\vee}(g)\lambda, v \rangle = \langle \lambda, \pi(g)v \rangle.$$



Constrained maximization: Statistical mechanics interpretation

Assume θ to be a kind of "microstate" with energy level $\ell(\theta)$. So Θ is some "state space".

Constrained maximization: Statistical mechanics interpretation

Assume θ to be a kind of "microstate" with energy level $\ell(\theta)$. So Θ is some "state space".

Statistical mechanics: Assume a distribution of the microstates (across "particles") maximizing entropy, constrained to have expected energy $\leq E_0$.

Constrained maximization: Statistical mechanics interpretation

Assume θ to be a kind of "microstate" with energy level $\ell(\theta)$. So Θ is some "state space".

Statistical mechanics: Assume a distribution of the microstates (across "particles") maximizing entropy, constrained to have expected energy $\leq E_0$.

Lagrange multiplier $\beta \geq 0$:

$$\underset{q \in \mathcal{P}_{\nu}(\Theta)}{\operatorname{arg min}} - \mathcal{H}_{\nu}(q) + \beta (\mathbb{E}_{q d \nu}[\ell] - E_0) = \underset{q \in \mathcal{P}_{\nu}(\Theta)}{\operatorname{arg min}} \mathbb{E}_{q d \nu}[\ell] - \frac{1}{\beta} \mathcal{H}_{\nu}(q)$$

 $au=rac{1}{eta}$ corresponds to the thermodynamical notion of temperature.



$$\arg\min_{q \in \mathcal{Q}} \mathbb{E}_{q d\nu}[\ell] - \tau \mathcal{H}(q) =$$

$$\underset{q \in \mathcal{Q}}{\arg \min} \, \mathbb{E}_{q d \nu}[\ell] - \tau \mathcal{H}(q) = \underset{q \in \mathcal{Q}}{\arg \min} \, \mathbb{E}_{q d \nu} \left[-\log e^{-\frac{1}{\tau}\ell} \right] + \mathbb{E}_{q d \nu}[\log q]$$

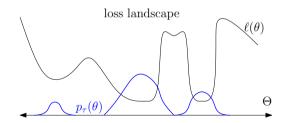
$$\underset{q \in \mathcal{Q}}{\arg\min} \, \mathbb{E}_{q d \nu}[\ell] - \tau \mathcal{H}(q) = \underset{q \in \mathcal{Q}}{\arg\min} \int_{\Theta} \log \left(\frac{q(\boldsymbol{\theta})}{e^{-\frac{1}{\tau}\ell(\boldsymbol{\theta})}} \right) q(\boldsymbol{\theta}) d\nu(\boldsymbol{\theta})$$

$$\mathop{\arg\min}_{q\in\mathcal{Q}}\mathbb{E}_{q\mathrm{d}\nu}[\ell] - \tau\mathcal{H}(q) = \mathop{\arg\min}_{q\in\mathcal{Q}}\int_{\Theta}\log\left(\frac{q(\boldsymbol{\theta})}{p_{\tau}(\boldsymbol{\theta})}\right)q(\boldsymbol{\theta})\mathrm{d}\nu(\boldsymbol{\theta}) + \mathsf{const.}$$

$$\underset{q \in \mathcal{Q}}{\arg\min} \, \mathbb{E}_{q d \nu}[\ell] - \tau \mathcal{H}(q) = \underset{q \in \mathcal{Q}}{\arg\min} \, \mathbb{D}_{\nu}(q \| p_{\tau}).$$

If $Q = \mathcal{P}_{\nu}(\Theta)$ then there is a unique minimizer $p_{\tau}(\theta) \propto e^{-\frac{1}{\tau}\ell(\theta)}$:

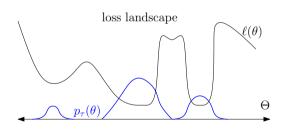
$$\underset{q \in \mathcal{Q}}{\arg\min} \, \mathbb{E}_{q d \nu}[\ell] - \tau \mathcal{H}(q) = \underset{q \in \mathcal{Q}}{\arg\min} \, \mathbb{D}_{\nu}(q \| p_{\tau}).$$



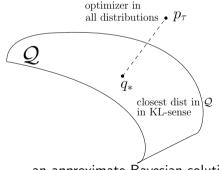
Minimize the objective $\mathcal{E}(q) := \mathbb{D}(q \| p_{\tau})$ for $q \in \mathcal{Q}$...

If $Q = \mathcal{P}_{\nu}(\Theta)$ then there is a unique minimizer $p_{\tau}(\theta) \propto e^{-\frac{1}{\tau}\ell(\theta)}$:

$$\arg\min_{q\in\mathcal{Q}} \mathbb{E}_{q\mathrm{d}\nu}[\ell] - \tau \mathcal{H}(q) = \arg\min_{q\in\mathcal{Q}} \mathbb{D}_{\nu}(q||p_{\tau}).$$



Minimize the objective $\mathcal{E}(q) := \mathbb{D}(q \| p_{\tau})$ for $q \in \mathcal{Q}...$



...an approximate Bayesian solution.

Let $\ell(\theta) = \sum_{i=1}^{N} \ell_i(\theta) + R(\theta)$. Observe new data $(\mathbf{x}_{\text{new}}, y_{\text{new}})$ with loss contribution ℓ_{new} .

Let $\ell(\theta) = \sum_{i=1}^{N} \ell_i(\theta) + R(\theta)$. Observe new data $(\mathbf{x}_{\text{new}}, y_{\text{new}})$ with loss contribution ℓ_{new} .

How to update p_{τ} ? Take $\tau = 1$

Let $\ell(\theta) = \sum_{i=1}^{N} \ell_i(\theta) + R(\theta)$. Observe new data $(\mathbf{x}_{\text{new}}, y_{\text{new}})$ with loss contribution ℓ_{new} .

How to update p_{τ} ? Take $\tau = 1$

Bayes' rule is about conditional probabilities, and updating priors:

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

Interpret $e^{-\ell_i(\theta)}$ as the likelihood of observing label y_i given the model parameter θ and \mathbf{x}_i . Interpret $\pi(\theta) \propto e^{-R(\theta)}$ as the prior on the parameters.

Let $\ell(\theta) = \sum_{i=1}^{N} \ell_i(\theta) + R(\theta)$. Observe new data $(\mathbf{x}_{\text{new}}, y_{\text{new}})$ with loss contribution ℓ_{new} .

How to update p_{τ} ? Take $\tau = 1$

Bayes' rule is about conditional probabilities, and updating priors:

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

Interpret $e^{-\ell_i(\theta)}$ as the likelihood of observing label y_i given the model parameter θ and \mathbf{x}_i . Interpret $\pi(\theta) \propto e^{-R(\theta)}$ as the prior on the parameters.

After one round of learning the posterior $p \propto e^{-\sum_i \ell_i \pi}$ is our prior belief about θ distribution. According to Bayes rule updated belief should be after a new data point.

$$p_{\mathsf{updated}}(\boldsymbol{\theta}) \propto e^{-\ell_{\mathsf{new}}(\boldsymbol{\theta})} p(\boldsymbol{\theta}).$$



Let $\ell(\theta) = \sum_{i=1}^{N} \ell_i(\theta) + R(\theta)$. Observe new data $(\mathbf{x}_{\text{new}}, y_{\text{new}})$ with loss contribution ℓ_{new} .

How to update p_{τ} ? Take $\tau = 1$

Bayes' rule is about conditional probabilities, and updating priors:

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

Interpret $e^{-\ell_i(\theta)}$ as the likelihood of observing label y_i given the model parameter θ and \mathbf{x}_i . Interpret $\pi(\theta) \propto e^{-R(\theta)}$ as the prior on the parameters.

After one round of learning the posterior $p \propto e^{-\sum_i \ell_i \pi}$ is our prior belief about θ distribution. According to Bayes rule updated belief should be after a new data point.

$$p_{\mathrm{updated}}(\boldsymbol{\theta}) \propto e^{-\ell_{\mathrm{new}}(\boldsymbol{\theta})} p(\boldsymbol{\theta}).$$

This is also the optimizer if we had initially considered the loss function $\ell_{updated} = \ell + \ell_{new}$.

Exponential Families

Let $T:\Theta\to V$, called the sufficient statistic. Call

$$\Omega = \Omega_{\nu}(T) = \left\{ \lambda \in V^{\vee} : A(\lambda) := \log \int_{\Theta} e^{-\langle \lambda, T(\boldsymbol{\theta}) \rangle} d\nu(\boldsymbol{\theta}) < \infty \right\}.$$

Then $q_{\lambda}(\boldsymbol{\theta}) = e^{-\langle \lambda, T(\boldsymbol{\theta}) \rangle - A(\lambda)}$ form an exponential family of distributions.

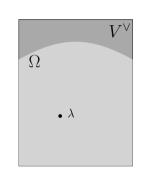
Exponential Families

Let $T:\Theta\to V$, called the sufficient statistic. Call

$$\Omega = \Omega_{\nu}(T) = \left\{ \lambda \in V^{\vee} : A(\lambda) := \log \int_{\Theta} e^{-\langle \lambda, T(\boldsymbol{\theta}) \rangle} d\nu(\boldsymbol{\theta}) < \infty \right\}.$$

Then $q_{\lambda}(\boldsymbol{\theta}) = e^{-\langle \lambda, T(\boldsymbol{\theta}) \rangle - A(\lambda)}$ form an exponential family of distributions.

$$\begin{split} -\frac{\partial A}{\partial \lambda_i} &= \int_{\Theta} T_i(\boldsymbol{\theta}) e^{-\langle \lambda, T(\boldsymbol{\theta}) \rangle - A(\lambda)} \mathrm{d}\nu(\boldsymbol{\theta}) = \mathbb{E}_{q_\lambda \mathrm{d}\nu}[T_i] =: \mu_i \\ \frac{\partial^2 A}{\partial \lambda_i \partial \lambda_j} &= \int_{\Theta} (T_i(\boldsymbol{\theta}) - \mu_i) (T_j(\boldsymbol{\theta}) - \mu_j) q_\lambda(\boldsymbol{\theta}) \mathrm{d}\nu(\boldsymbol{\theta}) \\ &= \mathbb{E}_{q_\lambda} \left[\left(\frac{\partial}{\partial \lambda_i} \log q_\lambda \right) \left(\frac{\partial}{\partial \lambda_j} \log q_\lambda \right) \right] =: F_{i,j}(\lambda) \quad \text{Fisher Matrix} \end{split}$$





Exponential Families

Let $T:\Theta\to V$, called the sufficient statistic. Call

$$\Omega = \Omega_{\nu}(T) = \left\{ \lambda \in V^{\vee} : A(\lambda) := \log \int_{\Theta} e^{-\langle \lambda, T(\boldsymbol{\theta}) \rangle} d\nu(\boldsymbol{\theta}) < \infty \right\}.$$

Then $q_{\lambda}(\boldsymbol{\theta}) = e^{-\langle \lambda, T(\boldsymbol{\theta}) \rangle - A(\lambda)}$ form an exponential family of distributions.

$$\begin{split} -\frac{\partial A}{\partial \lambda_i} &= \int_{\Theta} T_i(\boldsymbol{\theta}) e^{-\langle \lambda, T(\boldsymbol{\theta}) \rangle - A(\lambda)} \mathrm{d}\nu(\boldsymbol{\theta}) = \mathbb{E}_{q_\lambda \mathrm{d}\nu}[T_i] =: \mu_i \\ \frac{\partial^2 A}{\partial \lambda_i \partial \lambda_j} &= \int_{\Theta} (T_i(\boldsymbol{\theta}) - \mu_i) (T_j(\boldsymbol{\theta}) - \mu_j) q_\lambda(\boldsymbol{\theta}) \mathrm{d}\nu(\boldsymbol{\theta}) \\ &= \mathbb{E}_{q_\lambda} \left[\left(\frac{\partial}{\partial \lambda_i} \log q_\lambda \right) \left(\frac{\partial}{\partial \lambda_j} \log q_\lambda \right) \right] =: F_{i,j}(\lambda) \quad \text{Fisher Matrix} \end{split}$$

 V^{\vee} • λ

Example: If $T(\theta) = \begin{bmatrix} \theta \\ \theta^2 \end{bmatrix}$ then we get 1-D Gaussians $q_{\lambda}(\theta) \propto e^{-\lambda_1 \theta - \lambda_2 \theta^2}$ for $\lambda_2 > 0$.



What about exponential families on Θ which are closed under a Lie group action?

• Homogeneous space $\Theta \cong G/H$.

What about exponential families on Θ which are closed under a Lie group action?

- Homogeneous space $\Theta \cong G/H$.
- ν a relatively invariant base measure $\mathrm{d}\nu(g\cdot\theta)=\chi(g)\mathrm{d}\nu(\theta).$

What about exponential families on Θ which are closed under a Lie group action?

- Homogeneous space $\Theta \cong G/H$.
- ν a relatively invariant base measure $\mathrm{d}\nu(g\cdot\theta)=\chi(g)\mathrm{d}\nu(\theta).$
- A finite dimensional representation $\pi:G \to \operatorname{GL}(V)$.

What about exponential families on Θ which are closed under a Lie group action?

- Homogeneous space $\Theta \cong G/H$.
- ν a relatively invariant base measure $\mathrm{d}\nu(g\cdot\theta)=\chi(g)\mathrm{d}\nu(\theta).$
- A finite dimensional representation $\pi:G\to \operatorname{GL}(V).$
- A 1-cocycle of π such that $\alpha\big|_H \equiv 0$. So $\alpha:G \to V$ satisfies

$$\alpha(gh) = \pi(g)\alpha(h) + \alpha(g) = \alpha(g).$$

Thus $\alpha:\Theta\to V$.



What about exponential families on Θ which are closed under a Lie group action?

- Homogeneous space $\Theta \cong G/H$.
- ν a relatively invariant base measure $d\nu(g \cdot \theta) = \chi(g) d\nu(\theta)$.
- A finite dimensional representation $\pi: G \to \operatorname{GL}(V)$.
- A 1-cocycle of π such that $\alpha\big|_H \equiv 0$. So $\alpha:G \to V$ satisfies

$$\alpha(gh) = \pi(g)\alpha(h) + \alpha(g) = \alpha(g).$$

Thus $\alpha:\Theta\to V$.

Let
$$\lambda \in \Omega_{\nu}(\alpha) \subseteq V^{\vee}$$
 i.e.,
$$A(\lambda) = \log \int_{\Theta} e^{-\langle \lambda, \alpha(\theta) \rangle} \mathrm{d}\nu(\theta) < \infty.$$

$$q_{\lambda}(\theta) d\nu(\theta) := e^{-\langle \lambda, \alpha(\theta) \rangle - A(\lambda)} d\nu(\theta)$$

What about exponential families on Θ which are closed under a Lie group action?

- $\bullet \ \ {\rm Homogeneous \ space} \ \Theta \cong G/H.$
- ν a relatively invariant base measure $d\nu(g \cdot \theta) = \chi(g) d\nu(\theta)$.
- A finite dimensional representation $\pi:G\to \operatorname{GL}(V).$
- A 1-cocycle of π such that $\alpha\big|_H \equiv 0$. So $\alpha:G \to V$ satisfies

$$\alpha(gh) = \pi(g)\alpha(h) + \alpha(g) = \alpha(g).$$

Thus $\alpha:\Theta\to V$.

Let
$$\lambda \in \Omega_{\nu}(\alpha) \subseteq V^{\vee}$$
 i.e., $A(\lambda) = \log \int_{\Theta} e^{-\langle \lambda, \alpha(\theta) \rangle} d\nu(\theta) < \infty$.

$$q_{\lambda}(\theta) d\nu(\theta) := e^{-\langle \lambda, \alpha(\theta) \rangle - A(\lambda)} d\nu(\theta)$$



What about exponential families on Θ which are closed under a Lie group action?

- Homogeneous space $\Theta \cong G/H$.
- ν a relatively invariant base measure $d\nu(g \cdot \theta) = \chi(g) d\nu(\theta)$.
- A finite dimensional representation $\pi: G \to \operatorname{GL}(V)$.
- A 1-cocycle of π such that $\alpha\big|_H \equiv 0$. So $\alpha:G \to V$ satisfies

$$\alpha(gh) = \pi(g)\alpha(h) + \alpha(g) = \alpha(g).$$

Thus $\alpha:\Theta\to V$.

Let
$$\lambda \in \Omega_{\nu}(\alpha) \subseteq V^{\vee}$$
 i.e., $A(\lambda) = \log \int_{\Theta} e^{-\langle \lambda, \alpha(\theta) \rangle} d\nu(\theta) < \infty$.

$$q_{\lambda}(\theta) d\nu(\theta) := e^{-\langle \lambda, \alpha(\theta) \rangle - A(\lambda)} d\nu(\theta)$$

$$(q_{\lambda})^g(\theta) = \frac{1}{\chi(g)} q_{\lambda}(g^{-1}\theta)$$

What about exponential families on Θ which are closed under a Lie group action?

- Homogeneous space $\Theta \cong G/H$.
- ν a relatively invariant base measure $d\nu(g \cdot \theta) = \chi(g) d\nu(\theta)$.
- A finite dimensional representation $\pi: G \to \operatorname{GL}(V)$.
- A 1-cocycle of π such that $\alpha\big|_H \equiv 0$. So $\alpha:G \to V$ satisfies

$$\alpha(gh) = \pi(g)\alpha(h) + \alpha(g) = \alpha(g).$$

Thus $\alpha:\Theta\to V$.

Let
$$\lambda \in \Omega_{\nu}(\alpha) \subseteq V^{\vee}$$
 i.e., $A(\lambda) = \log \int_{\Omega} e^{-\langle \lambda, \alpha(\theta) \rangle} d\nu(\theta) < \infty$.

$$q_{\lambda}(\theta) d\nu(\theta) := e^{-\langle \lambda, \alpha(\theta) \rangle - A(\lambda)} d\nu(\theta)$$

$$(q_{\lambda})^g(\theta) = \frac{1}{\chi(g)} q_{\lambda}(g^{-1}\theta) = e^{-\langle \lambda, \alpha(g^{-1}\theta) \rangle - A(\lambda)}$$

What about exponential families on Θ which are closed under a Lie group action?

- Homogeneous space $\Theta \cong G/H$.
- ν a relatively invariant base measure $d\nu(g \cdot \theta) = \chi(g) d\nu(\theta)$.
- A finite dimensional representation $\pi: G \to \operatorname{GL}(V)$.
- A 1-cocycle of π such that $\alpha\big|_H \equiv 0$. So $\alpha:G \to V$ satisfies

$$\alpha(gh) = \pi(g)\alpha(h) + \alpha(g) = \alpha(g).$$

Thus $\alpha:\Theta\to V$.

Let
$$\lambda \in \Omega_{\nu}(\alpha) \subseteq V^{\vee}$$
 i.e., $A(\lambda) = \log \int_{\Omega} e^{-\langle \lambda, \alpha(\theta) \rangle} d\nu(\theta) < \infty$.

$$q_{\lambda}(\theta) d\nu(\theta) := e^{-\langle \lambda, \alpha(\theta) \rangle - A(\lambda)} d\nu(\theta)$$

forms an exponential family closed under pushforwards

$$(q_{\lambda})^{g}(\theta) = \frac{1}{\chi(g)} q_{\lambda}(g^{-1}\theta) = e^{-\langle \lambda, \alpha(g^{-1}\theta) \rangle - A(\lambda)}$$
$$= e^{-\langle \lambda, \pi(g^{-1})\alpha(\theta) \rangle - A(\lambda) - \alpha(g^{-1})}$$

35 / 36

What about exponential families on Θ which are closed under a Lie group action?

- Homogeneous space $\Theta \cong G/H$.
- ν a relatively invariant base measure $d\nu(g \cdot \theta) = \chi(g) d\nu(\theta)$.
- A finite dimensional representation $\pi: G \to \operatorname{GL}(V)$.
- A 1-cocycle of π such that $\alpha\big|_H \equiv 0$. So $\alpha:G \to V$ satisfies

$$\alpha(gh) = \pi(g)\alpha(h) + \alpha(g) = \alpha(g).$$

Thus $\alpha:\Theta\to V$.

Let
$$\lambda \in \Omega_{\nu}(\alpha) \subseteq V^{\vee}$$
 i.e., $A(\lambda) = \log \int_{\Theta} e^{-\langle \lambda, \alpha(\theta) \rangle} d\nu(\theta) < \infty$.

$$q_{\lambda}(\theta) d\nu(\theta) := e^{-\langle \lambda, \alpha(\theta) \rangle - A(\lambda)} d\nu(\theta)$$

$$(q_{\lambda})^{g}(\theta) = \frac{1}{\chi(g)} q_{\lambda}(g^{-1}\theta) = e^{-\langle \lambda, \alpha(g^{-1}\theta) \rangle - A(\lambda)}$$
$$= e^{-\langle \lambda, \pi(g^{-1})\alpha(\theta) \rangle - A(\lambda) - \alpha(g^{-1})}$$
$$\propto e^{-\langle \pi^{\vee}(g)\lambda, \alpha(\theta) \rangle}$$



What about exponential families on Θ which are closed under a Lie group action?

- Homogeneous space $\Theta \cong G/H$.
- ν a relatively invariant base measure $d\nu(g \cdot \theta) = \chi(g) d\nu(\theta)$.
- A finite dimensional representation $\pi: G \to \operatorname{GL}(V)$.
- A 1-cocycle of π such that $\alpha\big|_H \equiv 0$. So $\alpha:G \to V$ satisfies

$$\alpha(gh) = \pi(g)\alpha(h) + \alpha(g) = \alpha(g).$$

Thus $\alpha:\Theta\to V$.

Let
$$\lambda \in \Omega_{\nu}(\alpha) \subseteq V^{\vee}$$
 i.e., $A(\lambda) = \log \int_{\Theta} e^{-\langle \lambda, \alpha(\theta) \rangle} d\nu(\theta) < \infty$.

$$q_{\lambda}(\theta) d\nu(\theta) := e^{-\langle \lambda, \alpha(\theta) \rangle - A(\lambda)} d\nu(\theta)$$

$$(q_{\lambda})^{g}(\theta) = \frac{1}{\chi(g)} q_{\lambda}(g^{-1}\theta) = e^{-\langle \lambda, \alpha(g^{-1}\theta) \rangle - A(\lambda)}$$
$$= e^{-\langle \lambda, \pi(g^{-1})\alpha(\theta) \rangle - A(\lambda) - \alpha(g^{-1})}$$
$$\propto e^{-\langle \pi^{\vee}(g)\lambda, \alpha(\theta) \rangle} \propto q_{\pi^{\vee}(g)\lambda}(\theta)$$



The Lie group rule is for transformation families. BLR is NGD on λ 's of exponential families.

The Lie group rule is for transformation families. BLR is NGD on λ 's of exponential families. The overlap is harmonic exponential families:

Pushforwards of q_{λ} are still in the family,

$$(q_\lambda)^g = q_{\lambda'} \quad \text{ with } \quad \lambda' = \pi^\vee(g)\lambda$$

The Lie group rule is for transformation families. BLR is NGD on λ 's of exponential families. The overlap is harmonic exponential families:

Pushforwards of q_{λ} are still in the family,

$$(q_{\lambda})^g = q_{\lambda'} \quad \text{ with } \quad \lambda' = \pi^{\vee}(g)\lambda$$

Thus $\widetilde{\mathcal{Q}} = \{\lambda \in \Omega : \lambda = \pi^{\vee}(g)\lambda_0, g \in G\}$ and updates are given by $\lambda^{\text{updated}} = \pi^{\vee}(g^{\text{updated}})\lambda_0$.

The Lie group rule is for transformation families. BLR is NGD on λ 's of exponential families.

The overlap is harmonic exponential families:

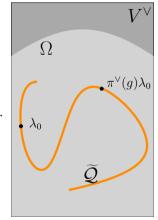
Pushforwards of q_{λ} are still in the family,

Other quantities of q_{λ} also vary with q:

$$(q_{\lambda})^g = q_{\lambda'}$$
 with $\lambda' = \pi^{\vee}(g)\lambda$

Thus $\widetilde{\mathcal{Q}} = \{\lambda \in \Omega : \lambda = \pi^{\vee}(g)\lambda_0, g \in G\}$ and updates are given by $\lambda^{\text{updated}} = \pi^{\vee}(g^{\text{updated}})\lambda_0$.

$$\mu(\lambda') = \pi(g)\lambda + \alpha(g)$$
$$A(\lambda') = A(\lambda) + \log(\chi(g)) + \alpha(g^{-1})$$



The Lie group rule is for transformation families. BLR is NGD on λ 's of exponential families.

The overlap is harmonic exponential families:

Pushforwards of q_{λ} are still in the family,

Other quantities of q_{λ} also vary with q:

$$(q_{\lambda})^g = q_{\lambda'}$$
 with $\lambda' = \pi^{\vee}(g)\lambda$

Thus $\widetilde{\mathcal{Q}} = \{\lambda \in \Omega : \lambda = \pi^{\vee}(g)\lambda_0, g \in G\}$ and updates are given by $\lambda^{\text{updated}} = \pi^{\vee}(g^{\text{updated}})\lambda_0$.

$$\mu(\lambda') = \pi(g)\lambda + \alpha(g)$$
$$A(\lambda') = A(\lambda) + \log(\chi(g)) + \alpha(g^{-1})$$

