

# PROBABILISTIC INFERENCE AND LEARNING

## LECTURE 11

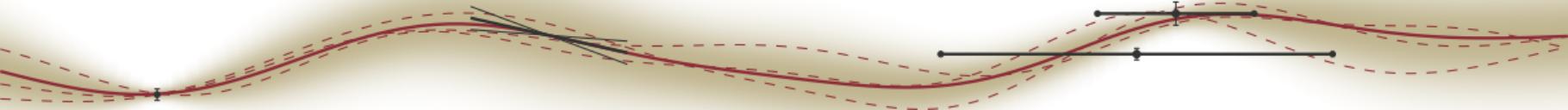
### EXPONENTIAL FAMILIES

Philipp Hennig

21 November 2018

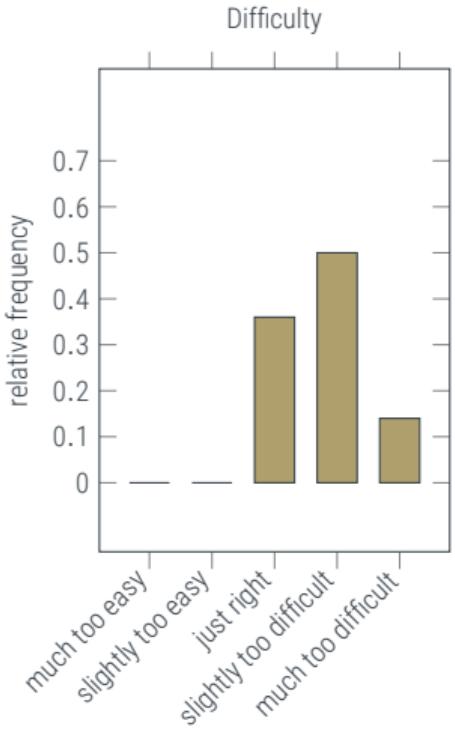
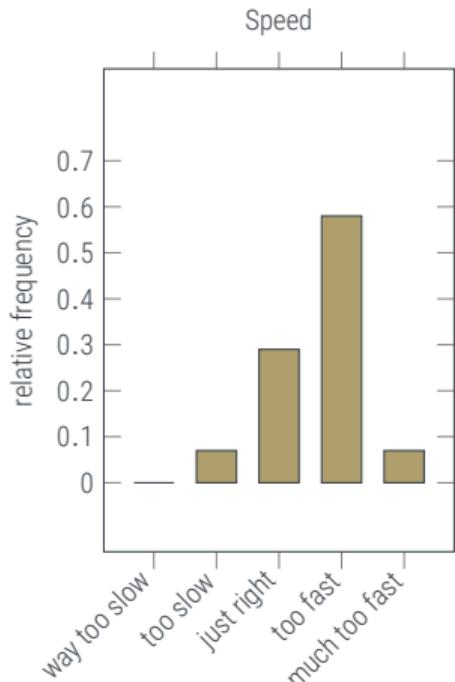
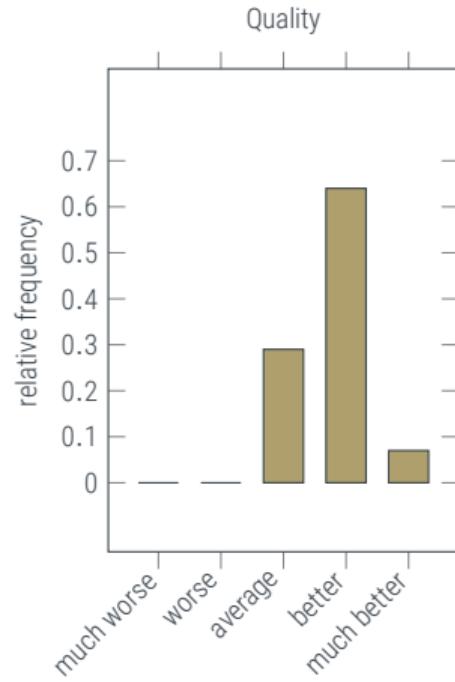


FACULTY OF SCIENCE  
DEPARTMENT OF COMPUTER SCIENCE  
CHAIR FOR THE METHODS OF MACHINE LEARNING



# Last Lecture: Debrief

Feedback dashboard





# Last Lecture: Debrief

## Detailed Feedback

### Things you did not like:

- ♦ the exercises are too hard

### Things you did not understand:

- ♦ steady-state
- ♦ what kind of object is an SDE?
- ♦ too many new concepts in one lecture

### Things you enjoyed:

- ♦ room for discussion
- ♦ recap
- ♦ the exercises are of a good difficulty level



## Overview of Lectures so far:

0. Introduction to Reasoning under Uncertainty
1. Probabilistic Reasoning
2. Probabilities over Continuous Variables
3. Gaussian Probability Distributions
4. Gaussian Parametric Regression
5. More on Parametric Regression
6. Gaussian Processes
7. More on Kernels & GPs
8. A practical GP example
9. Markov Chains, Time Series, Filtering
10. Classification
11. Empirical Example of Classification
12. Bayesianism and Frequentism
13. Stochastic Differential Equations
14. Exponential Families

## Today:

- A Toolbox of Probability Distributions
- A clean framework for efficient probabilistic inference, with some caveats



# Why is this hard?

The computational challenge in Bayesian Inference

$$p(x \mid y) = \frac{p(y \mid x)p(x)}{\int p(y \mid x)p(x) dx}$$



# Why is this hard?

The computational challenge in Bayesian Inference

$$p(x \mid y) = \frac{p(y \mid x)p(x)}{\int p(y \mid x)p(x) dx}$$

- the integral  $\int p(y \mid x)p(x) dx$  may be intractable
- thus, also expectations  $\int f(x)p(x \mid y) dx$  are hard
- even if we just want  $\arg \max_x p(x \mid y)$ , the optimization may be hard

# Hierarchical Bayesian Inference

Catch-up from previous lectures

Recall from GP regression: How to set parameters  $\theta$ ? From marginal likelihood  $p(Y | \theta)$ :

$$\begin{aligned}
 \hat{\boldsymbol{\theta}} &= \arg \max_{\boldsymbol{\theta}} \mathcal{N}(y; \phi_X^{\boldsymbol{\theta}^\top} \mu + b, \phi_X^{\boldsymbol{\theta}^\top} \Sigma \phi_X^{\boldsymbol{\theta}} + \Lambda) \\
 &= \arg \max_{\boldsymbol{\theta}} \log \mathcal{N}(y; \phi_X^{\boldsymbol{\theta}^\top} \mu + b, \phi_X^{\boldsymbol{\theta}^\top} \Sigma \phi_X^{\boldsymbol{\theta}} + \Lambda) \\
 &= \arg \min_{\boldsymbol{\theta}} -\log \mathcal{N}(y; \phi_X^{\boldsymbol{\theta}^\top} \mu + b, \phi_X^{\boldsymbol{\theta}^\top} \Sigma \phi_X^{\boldsymbol{\theta}} + \Lambda) \\
 &= \arg \min_{\boldsymbol{\theta}} \frac{1}{2} \left( \underbrace{(y - \phi_X^{\boldsymbol{\theta}^\top} \mu)^\top (\phi_X^{\boldsymbol{\theta}^\top} \Sigma \phi_X^{\boldsymbol{\theta}} + \Lambda)^{-1} (y - \phi_X^{\boldsymbol{\theta}^\top} \mu)}_{\text{square error}} + \underbrace{\log |\phi_X^{\boldsymbol{\theta}^\top} \Sigma \phi_X^{\boldsymbol{\theta}} + \Lambda|}_{\text{model complexity / Occam factor}} \right) + \frac{N}{2} \log 2\pi
 \end{aligned}$$

In general, hierarchical inference is not analytically tractable. However, there are special cases...



# Analytic Hierarchical Bayesian Inference

Inferring the Mean of a Gaussian

$$p(x | \mu) = \prod_{i=1}^n \mathcal{N}(x_i; \mu, \Sigma)$$

$$p(\mu | \mu_0, \Sigma_0) = \mathcal{N}(\mu; \mu_0, \Sigma_0)$$

$$\begin{aligned} p(\mu | x) &= \frac{p(x | \mu)p(\mu | \mu_0, \Sigma_0)}{p(x)} \\ &= \mathcal{N}(\mu; (\Sigma_0^{-1} + n\Sigma^{-1})^{-1}(\Sigma_0^{-1}\mu_0 + \Sigma^{-1} \sum_i x_i), (\Sigma_0^{-1} + n\Sigma^{-1})^{-1}) \end{aligned}$$

# Analytic Hierarchical Bayesian Inference

Inferring a Binary Distribution



$$p(x | f) = \prod_{i=1}^n f^x \cdot (1-f)^{1-x} \quad x \in \{0; 1\}$$
$$= f^{n_1} \cdot (1-f)^{n_0} \quad n_0 := n - n_1$$

$$p(f | \alpha, \beta) = \mathcal{B}(\alpha, \beta) = \frac{1}{B(\alpha, \beta)} f^{\alpha-1} (1-f)^{\beta-1}$$

$$p(f | x) = \mathcal{B}(\alpha + n_1, \beta + n_0) = \frac{1}{B(\alpha + n_1, \beta + n_0)} f^{\alpha+n_1-1} (1-f)^{\beta+n_0-1}$$



Pierre Simon, marquis de Laplace, 1749–1827

# Analytic Hierarchical Bayesian Inference

Inferring a Categorical Distribution



image: Deutsches Museum München

$$p(x) = \prod_{i=1}^n f_{x_i} \quad x \in \{0; \dots, K\}$$

$$= \prod_{k=1}^K f_k^{n_k} \quad n_k := |\{x_i \mid x_i = k\}|$$

$$p(f \mid \alpha) = \mathcal{D}(\alpha) = \frac{1}{B(\alpha)} \prod_{k=1}^K f_k^{\alpha_k - 1}$$

$$p(f \mid x) = \mathcal{D}(\alpha + n)$$



Peter Gustav Lejeune Dirichlet  
(1805–1859)



# Analytic Hierarchical Bayesian Inference

Inferring the (Co-) Variance of a Gaussian

$$p(\mathbf{x} \mid \sigma) = \prod_{i=1}^n \mathcal{N}(x_i; \mu, \sigma^2)$$
$$p(\sigma) = ?$$



# Analytic Hierarchical Bayesian Inference

Inferring the (Co-) Variance of a Gaussian

$$p(\mathbf{x} | \sigma) = \prod_{i=1}^n \mathcal{N}(x_i; \mu, \sigma^2)$$

$$p(\sigma) = ?$$

$$\log p(\mathbf{x} | \sigma) = -\frac{1}{2} \log \sigma^2 - \frac{1}{2} (\mathbf{x} - \mu)^2 \cdot \frac{1}{\sigma^2} - \frac{1}{2} \log 2\pi$$

# Analytic Hierarchical Bayesian Inference

Inferring the (Co-) Variance of a Gaussian

$$p(\mathbf{x} | \sigma) = \prod_{i=1}^n \mathcal{N}(x_i; \mu, \sigma^2)$$

$$p(\sigma) = ?$$

$$\log p(\mathbf{x} | \sigma) = -\frac{1}{2} \log \sigma^2 - \frac{1}{2} (\mathbf{x} - \mu)^T (\mathbf{x} - \mu) \cdot \frac{1}{\sigma^2} - \frac{1}{2} \log 2\pi$$

$$\log p(\sigma | \alpha, \beta) = (\alpha + 1) \log \sigma^{-2} - \beta \cdot \frac{1}{\sigma^2} - Z(\alpha, \beta)$$

$$p(\sigma | \alpha, \beta) = \frac{\beta^\alpha}{\Gamma(\alpha)} (\sigma^{-2})^{\alpha+1} e^{-\beta\sigma^{-2}} =: \mathcal{G}(\sigma^{-2}; \alpha, \beta)$$

$$p(\sigma | \alpha, \beta, \mathbf{x}) = \mathcal{G}\left(\sigma^{-2}; \alpha + \frac{n}{2}, \beta + \frac{1}{2} \sum_i (x_i - \mu)^2\right)$$



Daniel Bernoulli (1700–1782)

# Analytic Hierarchical Bayesian Inference

Inferring Mean and Co-Variance of a Gaussian

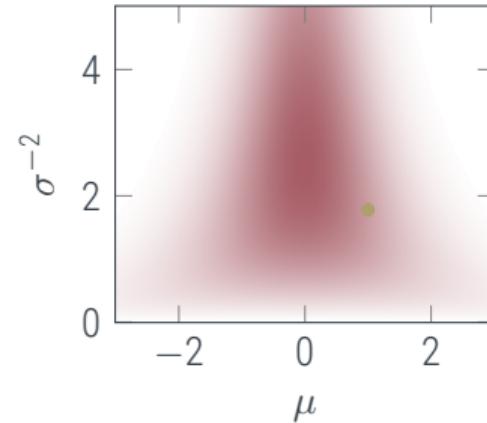
$$p(\mathbf{x} | \mu, \sigma) = \prod_{i=1}^n \mathcal{N}(x_i; \mu, \sigma^2)$$

$$p(\mu, \sigma | \mu_0, \nu, \alpha, \beta) = \mathcal{N}\left(\mu; \mu_0, \frac{\sigma^2}{\nu}\right) \mathcal{G}(\sigma^{-2}; \alpha, \beta)$$

$$p(\mu, \sigma | \mathbf{x}, \mu_0, \nu, \alpha, \beta) = \mathcal{N}\left(\mu; \frac{\nu\mu_0 + n\bar{x}}{\nu + n}, \frac{\sigma^2}{\nu + n}\right).$$

$$\mathcal{G}\left(\sigma^{-2}; \alpha + \frac{n}{2}, \beta + \frac{1}{2} \sum_{i=1}^n (x_i - \bar{x})^2 + \frac{n\nu}{2(n+\nu)} (\bar{x} - \mu_0)^2\right)$$

$$\text{where } \bar{x} := \frac{1}{n} \sum_{i=1}^n x_i$$





# Analytic Hierarchical Bayesian Inference

Inferring Mean and Co-Variance of a Gaussian

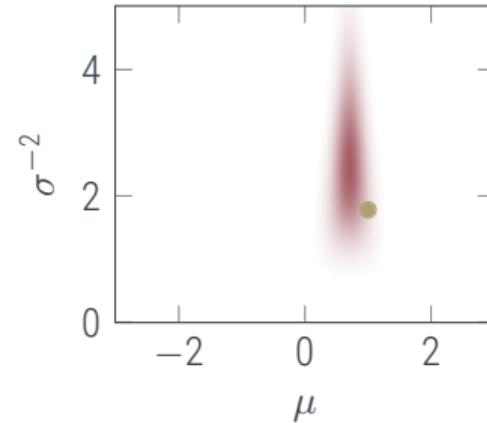
$$p(\mathbf{x} | \mu, \sigma) = \prod_{i=1}^n \mathcal{N}(x_i; \mu, \sigma^2)$$

$$p(\mu, \sigma | \mu_0, \nu, \alpha, \beta) = \mathcal{N}\left(\mu; \mu_0, \frac{\sigma^2}{\nu}\right) \mathcal{G}(\sigma^{-2}; \alpha, \beta)$$

$$p(\mu, \sigma | \mathbf{x}, \mu_0, \nu, \alpha, \beta) = \mathcal{N}\left(\mu; \frac{\nu\mu_0 + n\bar{x}}{\nu + n}, \frac{\sigma^2}{\nu + n}\right).$$

$$\mathcal{G}\left(\sigma^{-2}; \alpha + \frac{n}{2}, \beta + \frac{1}{2} \sum_{i=1}^n (x_i - \bar{x})^2 + \frac{n\nu}{2(n+\nu)} (\bar{x} - \mu_0)^2\right)$$

$$\text{where } \bar{x} := \frac{1}{n} \sum_{i=1}^n x_i$$





# Conjugate Prior Inference

a beautiful idea, not to be underestimated

## Definition (Conjugate Prior)

Let  $D$  and  $x$  be a data-set and a variable to be inferred, respectively, connected by the likelihood  $p(D | x) = \ell(D; x)$ . A **conjugate prior to  $\ell$  for  $x$**  is a probability measure with pdf  $p(x) = \pi(x; \theta)$  of functional form  $\pi$ , such that

$$p(x | D) = \frac{\ell(D; x)\pi(x; \theta)}{\int \ell(D; x)\pi(x; \theta) dx} = \pi(x; \theta').$$

That is, such that the posterior arising from  $\ell$  is of the same functional form as the prior, with updated parameters.

- E. Pitman. *Sufficient statistics and intrinsic accuracy* (1936). Math. Proc. Cambr. Phil. Soc. 32(4), 1936.  
P. Diaconis and D. Ylvisaker, *Conjugate priors for exponential families*. Annals of Statistics 7(2), 1979.



- **Conjugate priors** allow analytic Bayesian inference
- How can we construct them in general?



# Exponential Families

## Exponentials of a Linear Form

### Definition (Exponential Family, simplified form)

Consider a random variable  $X$  taking values  $x \in \mathbb{X} \subset \mathbb{R}^n$ . A probability distribution for  $X$  with pdf of the functional form

$$p_w(x) = \exp [\phi(x)^T w - \log Z(w)] = \frac{1}{Z(w)} e^{\phi(x)^T w} = p(x | w)$$

is called an **exponential family** of probability measures. The function  $\phi : \mathbb{X} \rightarrow \mathbb{R}^d$  is called the **sufficient statistics**. The parameters  $w \in \mathbb{R}^d$  are the **natural parameters** of  $p_w$ . The normalization constant  $Z(w) : \mathbb{R}^d \rightarrow \mathbb{R}$  is the **partition function**.



# Exponential Families

## Exponentials of a Linear Form

[Diaconis & Ylvisaker, AoS 7/2, 1979]

### Definition (Exponential Family, formal definition, only if you care for it)

Consider a  $\sigma$ -finite measure  $\mu$  on the Borel sets of  $\mathbb{R}^d$ . Consider the convex hull of the support set of  $\mu$ , and let  $\mathbb{X}$  be the interior of this convex set, assumed to be nonempty and open in  $\mathbb{R}^d$ . For  $w \in \mathbb{R}^d$ , define  $Z(w) = \ln \int e^{w \cdot x} d\mu(x)$  and let  $W = \{w \mid Z(w) < \infty\}$ . Hölder's inequality shows that  $\Theta$  is a convex set, called the **natural parameter** space. Assume further that  $W$  is nonempty and open in  $\mathbb{R}^d$ . The **exponential family**  $\{P_w\}$  of probability measures **through**  $\mu$  is determined by

$$dP_w(x) = e^{x \cdot w - Z(w)} d\mu(x).$$

Nb: Relative to our more practically-minded definition on the previous slide, this definition absorbs the sufficient statistics  $\phi$  into the definition of  $X$ .

# A Family Meeting

incomplete list of exponential families



Bernoulli	$\phi(x) = [x]$	$\mathbb{X} = \{0; 1\}$
Poisson	$\phi(x) = [x]$	$\mathbb{X} = \mathbb{R}_+$
Laplace	$\phi(x) = [1, x]^\top$	$\mathbb{X} = \mathbb{R}$
$\chi^2$	$\phi(x) = [x, -\log x]$	$\mathbb{X} = \mathbb{R}$
Dirichlet	$\phi(x) = [\log x]$	$\mathbb{X} = \mathbb{R}_+$
Euler ( $\Gamma$ )	$\phi(x) = [x, \log x]$	$\mathbb{X} = \mathbb{R}_+$
Wishart	$\phi(X) = [X, \log  X ]$	$\mathbb{X} = \{X \in \mathbb{R}^{N \times N} \mid v^\top X v \geq 0 \forall v \in \mathbb{R}^N\}$
Gauss	$\phi(X) = [X, XX^\top]$	$\mathbb{X} = \mathbb{R}^N$
Boltzmann	$\phi(X) = [X, \text{triag}(XX^\top)]$	$\mathbb{X} = \{0; 1\}^N$



# Exponential Families have Conjugate Priors

but the prior's normalization constant can be tricky

- Consider the exponential family  $p_w(x | w) = \exp [\phi(x)^\top w - \log Z(w)]$

- its conjugate prior is the exponential family

$$F(\alpha, \nu) = \int \exp(\alpha^\top w - \nu^\top \log Z(w)) dw$$

$$p_\alpha(w | \alpha, \nu) = \exp \left[ \begin{pmatrix} w \\ -\log Z(w) \end{pmatrix}^\top \begin{pmatrix} \alpha \\ \nu \end{pmatrix} - \log F(\alpha, \nu) \right]$$

$$\text{because } p_\alpha(w | \alpha, \nu) \prod_{i=1}^n p_w(x_i | w) \propto p_\alpha \left( w \middle| \alpha + \sum_i \phi(x_i), \nu + n \right)$$

- and the predictive is

$$\begin{aligned} p(x) &= \int p_w(x | w) p_\alpha(w | \alpha, \nu) dw = \int e^{(\phi(x) + \alpha)^\top w + (\nu + 1) \log Z(w) + \log F(\alpha, \nu)} dw \\ &= \frac{F(\phi(x) + \alpha, \nu + 1)}{F(\alpha, \nu)} \end{aligned}$$

Computing  $F(\alpha, \nu)$  can be tricky. In general, this is **the** challenge when constructing an EF.

# Example: The Gaussian distribution

There it is, again!

$$\begin{aligned}
 p_w(x | w) &= \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right) = \mathcal{N}(x; \mu, \sigma^2) \\
 &= \exp\left(-\frac{x^2 - 2\mu x + \mu^2}{2\sigma^2} - \log \sqrt{2\pi\sigma^2}\right) \\
 &= \exp\left(\begin{bmatrix} x & -1/2 x^2 \end{bmatrix} \begin{bmatrix} \mu/\sigma^2 \\ 1/\sigma^2 \end{bmatrix} - \left(\frac{\mu^2}{2\sigma^2} + \log \sqrt{2\pi\sigma^2}\right)\right) \\
 &= \exp\left(\begin{bmatrix} \phi_1(x) & \phi_2(x) \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} - \underbrace{\left(\frac{w_1^2}{2w_2} - \frac{1}{2} \log w_2 + \log \sqrt{2\pi}\right)}_{\log Z(w)}\right)
 \end{aligned}$$

- The natural parameters are **precision**  $\sigma^{-2}$  and **precision-adjusted mean**  $\mu\sigma^{-2}$
- the sufficient statistics are the sample *mean* and  $(-1/2)$  *variance*
- The conjugate prior is the *Normal-Gamma*, the predictive marginal is the *Student-t* distribution



why they're called **sufficient statistics**

- Consider the exponential family

$$p_w(x | w) = \exp [\phi(x)^T w - \log Z(w)]$$

- for iid data:

$$p_w(x_1, x_2, \dots, x_n | w) = \prod_i^n p_w(x_i | w) = \exp \left( \sum_i^n \phi^T(x_i) w - n \log Z(w) \right)$$

- to find the **maximum likelihood estimate** for  $w$ , set

$$\nabla_w \log p(x | w) = 0 \quad \Rightarrow \quad \nabla_w \log Z(w) = \frac{1}{n} \sum_i \phi(x_i)$$

- hence, collect **statistics** of  $\phi$ , compute  $\nabla_w \log Z(w)$  and solve the above for  $w$ .

# Other great properties

exponential families make many things easy

- Re-phrased from above: because  $\int_{\mathbb{X}} dp_w(x) = 1$ , we have

$$\begin{aligned} \nabla_w \int p_w(x | w) dx &= \int \nabla_w p_w(x | w) dx &= \int \phi(x) dp_w(x | w) - \nabla_w \log Z(w) \int dp_w(x | w) \\ &= \nabla_w 1 &= 0 \\ \Rightarrow \mathbb{E}_{p_w}(\phi(x)) &= \nabla_w \log Z(w) \end{aligned}$$

- hence, if we should need to compute  $\mathbb{E}_{p_w}(\phi(x))$ , we can do so by *differentiating*  $\log Z$  wrt.  $w$  instead of *integrating*  $p$  over  $x$ . (actually, we're efficiently re-using someone else's integral)
- Note that an exponential family forms a *Abelian semigroup* on  $w$ :

$$p_w(x | w_1) \cdot p_w(x | w_2) \propto p(x | w_1 + w_2)$$

- Thus, combining information about  $x$  from independent  $p_w$ -sources can be done by floating point addition. In this sense, exponential families map inference to addition.



Can we use exponential families  $p_w(x) = e^{\phi(x)^T w} / Z(w)$  to learn **distributions**,  
just like we used linear forms  $f(x) = \phi(x)^T w$  to learn **functions**?

Yes! In fact, we can even do **Bayesian** distribution regression. It is called *conjugate prior inference*.

# Recap: Regression on Functions

The  $\ell_2$  loss

- Recall lectures 4–8: **regression on real functions**:

Given  $(y_i, x_i)_{i=1,\dots,n}$ , and assuming  $y_i = f(x_i)$ , assume  $f(x) = \phi(x)^\top w$  and find

$$\hat{f}(x) = \arg \min_{w \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n \|y_i - \phi(x_i)^\top w\|^2 + \|w\|_\Sigma^2 =: \arg \min_{w \in \mathbb{R}^d} \mathcal{L}_2(w)$$

- can interpret  $\exp -\mathcal{L}_2(w)$  as the (unnormalized) **posterior** on  $w$  from the Gaussian prior  $\exp(-\|w\|_\Sigma^2)$ .
- assume  $x_i \sim p(x)$ , then the Loss approximates an expected log posterior

$$\hat{f} = \arg \min_{w \in \mathbb{R}^d} \int \|f(x) - \phi(x)^\top w\|^2 dp(x) + \|w\|_\Sigma^2$$



# Interlude: KL divergence

The most mis-spelled names in statistics

## Definition (Kullback-Leibler divergence)

Let  $P$  and  $Q$  be probability distributions over  $\mathbb{X}$  with pdf's  $p(x)$  and  $q(x)$ , respectively. The **KL-divergence from  $Q$  to  $P$**  is defined as

$$D_{\text{KL}}(P||Q) := \int \log \left( \frac{p(x)}{q(x)} \right) dp(x)$$

(I will often write  $D_{\text{KL}}(p||q)$  instead)

Some properties:

- $D_{\text{KL}}(P||Q) \neq D_{\text{KL}}(Q||P)$
- $D_{\text{KL}}(P||Q) \geq 0, \forall P, Q$  (**Gibbs' inequality**), and
- $D_{\text{KL}}(P||Q) = 0 \Leftrightarrow p \equiv q$  almost everywhere



Solomon Kullback  
(1907–1994)



Richard Leibler  
(1914–2003)

# Regression on Distributions!

Fitting distributions with exponential families

- Given  $[x_i]_{i=1,\dots,n}$  with  $x_i \sim p(x)$ , assume

$$p(x) \approx \hat{p}(x \mid w) = \exp(\phi(x)^\top w - \log Z(w))$$

- to find  $\hat{w}$ , consider

$$\begin{aligned}\hat{w} &= \arg \min_{w \in \mathbb{R}^d} D_{\text{KL}}(p(x) \parallel \hat{p}(x \mid w)) = \arg \min_{w \in \mathbb{R}^d} \int [\log p(x) - \log \hat{p}(x \mid w)] dp(x) \\ &= \arg \min_{w \in \mathbb{R}^d} \underbrace{\int \log p(x) dp(x)}_{-\mathbb{H}(p)} + \mathbb{E}_p(\phi(x))^\top w + \log Z(w) = \arg \min_{w \in \mathbb{R}^d} \mathcal{L}_{\log}(w)\end{aligned}$$

- Find minimum at  $\nabla_w \mathcal{L}_{\log}(w) = 0$ , where

$$\frac{1}{n} \sum_{i=1}^n \phi(x_i) \approx \mathbb{E}_p(\phi(x)) = -\nabla_w \log Z(w) = \mathbb{E}_{\hat{p}}(\phi(x))$$

# Regression on Distributions!

Fitting distributions with exponential families

- Given  $[x_i]_{i=1,\dots,n}$  with  $x_i \sim p(x)$ , assume

$$p(x) \approx \hat{p}(x \mid w) = \exp(\phi(x)^\top w - \log Z(w))$$

- to find  $\hat{w}$ , consider (to regularize, include the conjugate prior. No need to know its normalizer!)

$$\begin{aligned}\hat{w} &= \arg \min_{w \in \mathbb{R}^d} D_{\text{KL}}(p(x) \parallel \hat{p}(x, w)) = \arg \min_{w \in \mathbb{R}^d} \int [\log p(x) - \log \hat{p}(x \mid w)] dp(x) + \alpha^\top w + \nu \log Z(w) \\ &= \arg \min_{w \in \mathbb{R}^d} \underbrace{\int \log p(x) dp(x)}_{-\mathbb{H}(p)} + \mathbb{E}_p(\phi(x))^\top w + \log Z(w) + \alpha^\top w + \nu \log Z(w) = \arg \min_{w \in \mathbb{R}^d} \tilde{\mathcal{L}}_{\log}(w)\end{aligned}$$

- Find minimum at  $\nabla_w \tilde{\mathcal{L}}_{\log}(w) = 0$ , where

$$\frac{1}{n} \sum_{i=1}^n \phi(x_i) \approx \mathbb{E}_p(\phi(x)) = -(1 + \nu) \nabla_w \log Z(w) - \alpha$$

- If the conjugate prior has tractable normalizer, we can do full Bayesian inference on  $w$ !

assuming  $x \sim p_w(x | w) = \exp [\phi(x)^T w - \log Z(w)]$

and assuming  $p_\alpha(w | \alpha, \nu) = \exp \left[ \begin{pmatrix} \alpha \\ \nu \end{pmatrix}^T \begin{pmatrix} w \\ \log Z(w) \end{pmatrix} - \log F(\alpha, \nu) \right]$

then  $p_\alpha(w | \alpha, \nu) \prod_{i=1}^n p_w(x_i | w) \propto p_\alpha \left( w \middle| \alpha + \sum_i \phi(x_i), \nu + n \right)$

- If the conjugate prior is intractable, we can still do MAP:

choosing  $w$  such that  $\frac{1}{n} \sum_{i=1}^n \phi(x_i) \approx \mathbb{E}_p(\phi(x)) = -(1 + \nu) \nabla_w \log Z(w) - \alpha$

equates to setting  $w = \arg \max_w \frac{p_w(x | w) p_\alpha(w | \alpha, \nu)}{\int p_w(x | w) p_\alpha(w | \alpha, \nu) dw}$



# Wouldn't you want to join this club?

Build your own exponential family!



# Building our own Exponential Family

just for fun

- choose features (come up with grand motivation:  
attraction/repulsion)

$$\phi(x) = \begin{bmatrix} -x^2 \\ -x^{-2} \end{bmatrix}$$

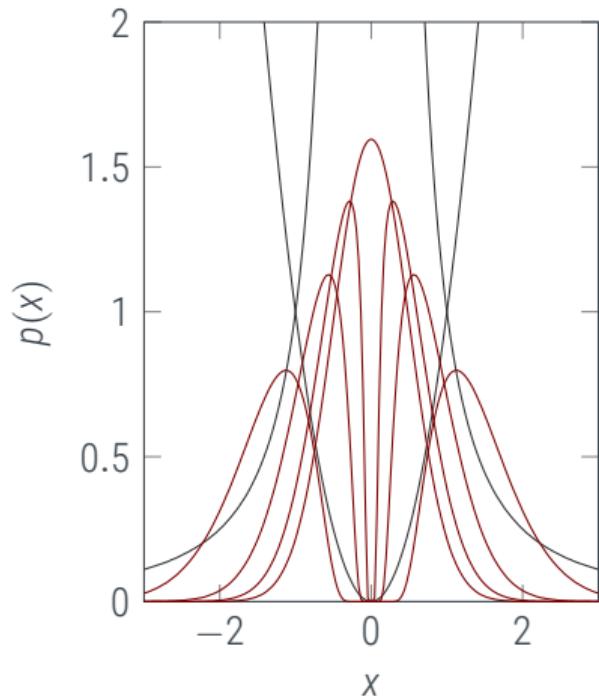
- solve integral (the hard bit)

$$Z(w) = \int_0^\infty \exp(-w_1 x^2 - w_2/x^2) dx = \sqrt{\frac{\pi}{w_1}} e^{-2\sqrt{w_1 w_2}}$$

- profit! The **bagel-distribution!**

$$\mathcal{H}(x; w) = \sqrt{\frac{w_1}{\pi}} e^{2\sqrt{w_1 w_2}} e^{-w_1 x^2 - w_2/x^2}$$

- don't know the conjugate prior, though. :(



# Let's fit a distribution!

collecting sufficient statistics

- We need

$$\log Z(w) = -2(w_1 w_2)^{1/2} - \frac{1}{2} \log w_1 + \frac{1}{2} \log \pi$$

$$-\nabla_w \log Z(w) = \begin{bmatrix} \sqrt{\frac{w_2}{w_1}} + \frac{1}{2w_1} \\ \sqrt{\frac{w_1}{w_2}} \end{bmatrix} \stackrel{!}{=} -\frac{1}{n} \sum_i \begin{bmatrix} x_i^2 \\ x_i^{-2} \end{bmatrix} =: \begin{bmatrix} \bar{\mu} \\ \bar{\omega} \end{bmatrix}$$

$$\Rightarrow \hat{w}_1 = \frac{1}{2(\bar{\mu} - \bar{\omega})} \quad \hat{w}_2 = \frac{\hat{w}_1}{\bar{\omega}^2}$$

## Summary:



- Conjugate Priors allow analytic inference of “nuisance parameters” in probabilistic models
- Exponential Families
  - guarantee the existence of conjugate priors, although not always tractable ones
  - allow analytic MAP inference from only a finite set of *sufficient statistics*

Conjugate prior inference with exponential families is a form of Bayesian **regression on distributions**. Gaussian process inference, in this sense, is inference on the unknown mean of a Gaussian distribution.

- The hardest part is finding the normalization constant. In fact, finding the normalization constant is *the only* hard part.
- Exponential families are a way to turn someone else’s integral into an inference algorithm!

## Summary:



- Conjugate Priors allow analytic inference of “nuisance parameters” in probabilistic models
- Exponential Families
  - guarantee the existence of conjugate priors, although not always tractable ones
  - allow analytic MAP inference from only a finite set of *sufficient statistics*

Conjugate prior inference with exponential families is a form of Bayesian **regression on distributions**. Gaussian process inference, in this sense, is inference on the unknown mean of a Gaussian distribution.

- The hardest part is finding the normalization constant. In fact, finding the normalization constant is *the only* hard part.
- Exponential families are a way to turn someone else’s integral into an inference algorithm!

If you happen to know (a good approximation to)

$$Z = \int_{-\infty}^{\infty} \exp \left[ \sum_i w_i \exp \left( -\frac{(x - c_i)^2}{2\lambda^2} \right) \right] dx$$

I may have a job for you ...