# Probabilistic Machine Learning
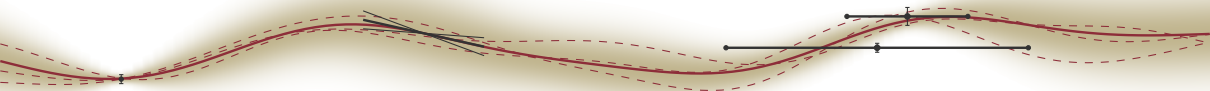## Lecture 05
## Exponential Families

Philipp Hennig

04 May 2023

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN

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

- ▶ "bonus points" means not everyone should get them
- ▶ otherwise, the exam just gets a steeper grading curve, and nobody gains anything
- ▶ *It's ok* to not be able to finish the exercise. You do not need all the bonus points to finish the exam!
- ▶ please do not post solutions on the forum. Sometimes they even add confusion!
- ▶ keep in mind that some things can actually be computed by hand. Auto-diff and optimization are powerful tools, but they do have computational and coding cost.

$$p(w \mid x) = \frac{p(x \mid w)p(w)}{\int p(x \mid w)p(w)\, dw} \overset{\text{"often"}}{=} \frac{\prod_{n=1}^{N} p(x_n \mid w)p(w)}{\int \prod_{n=1}^{N} p(x_n \mid w)p(w)\, dw}$$

This is very abstact. It would be nice to reduce this to

▶ get some **data** $x \in \mathbb{R}^{N \times D}$ ($N$: "batch dim", $D$: "input dim")

▶ compute some **statistics** $\phi(x) \in \mathbb{R}^{F}$ (note this can consume either or both axes of $x$) that capture the **algebraic structure** of $p(x \mid w)$

▶ *somehow* compute $p(w \mid x)$, using only $\phi(x)$ and nothing else about the data.

### Definition (Conjugate Prior)

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

$$p(w \mid x) = \frac{\ell(x; w)g(w; \theta)}{\int \ell(x; w)g(w; \theta)\, dw} = g(w; \theta + \phi(x)).$$

That is, such that the posterior arising from $\ell$ is of the same functional form as the prior, with updated parameters arising by **adding** some **sufficient statistics** of the observation $x$ to the prior's 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, by mapping it ot
  ► "data processing" $\phi(x) : \mathbb{R}^{N \times D} \to \mathbb{R}^F$
  ► "inference" $g(w; \theta + \phi(x))$, simple (floating-point) addition yields the full, normalized posterior!
  ► (the implicit assumption in this formulation is that $g$ can be computed in decent time)
► How can we construct them in general?

### 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

$$p_w(x) = h(x) \exp\left[\phi(x)^\mathsf{T} w - \log Z(w)\right] = \frac{h(x)}{Z(w)} e^{\phi(x)^\mathsf{T} w} = p(x \mid w)$$

is called an **exponential family** of probability measures. The function $\phi : \mathbb{X} \to \mathbb{R}^d$ is called the **sufficient statistics**. The parameters $w \in \mathbb{R}^d$ are the **natural parameters** of $p_w$. The normalization constant $Z(w) : \mathbb{R}^d \to \mathbb{R}$ is the **partition function**. The function $h(x) : \mathbb{X} \to \mathbb{R}_+$ is the **base measure**. For notational convenience, it can be useful to re-parametrize the natural parameters $w$ as $w := \eta(\theta)$ in terms of *canonical parameters* $\theta$.

Exponential families provide the probabilistic analogue to data types

| Name | sufficient stats | domain | use case |
|---|---|---|---|
| Bernoulli | $\phi(x) = [x]$ | $\mathbb{X} = \{0; 1\}$ | coin toss |
| Poisson | $\phi(x) = [x]$ | $\mathbb{X} = \mathbb{R}_+$ | emails per day |
| Laplace | $\phi(x) = [1, x]^{\mathsf{T}}$ | $\mathbb{X} = \mathbb{R}$ | floods |
| Helmert ($\chi^2$) | $\phi(x) = [x, -\log x]$ | $\mathbb{X} = \mathbb{R}$ | variances |
| Dirichlet | $\phi(x) = [\log x]$ | $\mathbb{X} = \mathbb{R}_+$ | class probabilities |
| Euler ($\Gamma$) | $\phi(x) = [x, \log x]$ | $\mathbb{X} = \mathbb{R}_+$ | variances |
| Wishart | $\phi(X) = [X, \log |X|]$ | $\mathbb{X} = \{X \in \mathbb{R}^{N \times N} \mid v^{\mathsf{T}} X v \geq 0 \forall v \in \mathbb{R}^N\}$ | covariances |
| Gauss | $\phi(X) = [X, XX^{\mathsf{T}}]$ | $\mathbb{X} = \mathbb{R}^N$ | functions |
| Boltzmann | $\phi(X) = [X, \text{triag}(XX^{\mathsf{T}})]$ | $\mathbb{X} = \{0; 1\}^N$ | thermodynamics |

Note: *Each row* of this table is *one* exponential family. Some authors (e.g. Murphy) call the entire table "the exponential family". We will *not* do this.

# CODE

– Thanks to Marvin Pförtner for pair coding –

▶ Consider the exponential family $p_w(x \mid w) = h(x) \exp[\phi(x)^\intercal w - \log Z(w)]$

▶ its conjugate prior is the exponential family $\qquad F(\alpha, \nu) := \int \exp(\alpha^\intercal w - \nu \log Z(w)) \, dw$

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

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

Computing $F(\alpha, \nu)$ can be tricky. But if we have it, it completely automates Bayesian inference!
Finding $F$ is thus the challenge when constructing an EF.

▶ Consider the exponential family $p_w(x \mid w) = h(x) \exp\left[\phi(x)^\mathsf{T} w - \log Z(w)\right]$

▶ its conjugate prior is the exponential family $\qquad F(\alpha, \nu) := \int \exp(\alpha^\mathsf{T} w - \nu \log Z(w))\, dw$

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

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

▶ and the predictive is

$$p(x) = \int p_w(x \mid w) p_\alpha(w \mid \alpha, \nu)\, dw = h(x) \int e^{(\phi(x)+\alpha)^\mathsf{T} w + (\nu+1)\log Z(w) - \log F(\alpha, \nu)}\, dw$$

$$= h(x) \frac{F(\phi(x) + \alpha, \nu + 1)}{F(\alpha, \nu)}$$

Computing $F(\alpha, \nu)$ can be tricky. But **if** we have it, it completely automates Bayesian inference!
Finding $F$ is thus **the** challenge when constructing an EF.

▶ Once we decide to use a particular **exponential family** $(h(x), \phi(x))$ as the **model** for some data $x$, we *automatically get* a *conjugate prior* in the form of another exponential family

▶ in principle, this provides the means to do analytic Bayesian inference, **if** we can compute the partition function $Z(w)$ of the likelihood and $F(\alpha, \nu)$, that of the conjugate prior.

▶ This solves the *learning* problem, i.e. how to extract information from the data. It does not guarantee that we will be able to compute moments, or sample from the posterior, but these can be encapsulated to subroutines, now that the data is dealt with.

CODE

$$p_w(x \mid w) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right) = \mathcal{N}(x; \mu, \sigma^2)$$

$$= \exp\left(-\frac{x^2 - 2\mu x + \mu^2}{2\sigma^2} - \log\sqrt{2\pi\sigma^2}\right)$$

$$= \exp\left(\begin{bmatrix} x & -1/2\, x^2 \end{bmatrix} \begin{bmatrix} \mu/\sigma^2 \\ 1/\sigma^2 \end{bmatrix} - \left(\frac{\mu^2}{2\sigma^2} + \log\sqrt{2\pi\sigma^2}\right)\right)$$

$$= \exp\left(\begin{bmatrix} \phi_1(x) & \phi_2(x) \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} - \underbrace{\frac{1}{2}\left(\frac{w_1^2}{w_2} - \log w_2 + \log(2\pi)\right)}_{\log Z(w)}\right)$$

▶ The natural parameters are **precision** $\sigma^{-2}$ and **precision-adjusted mean** $\mu\sigma^{-2}$
▶ the sufficient statistics are the *first two (non-central) sample moments* of the data
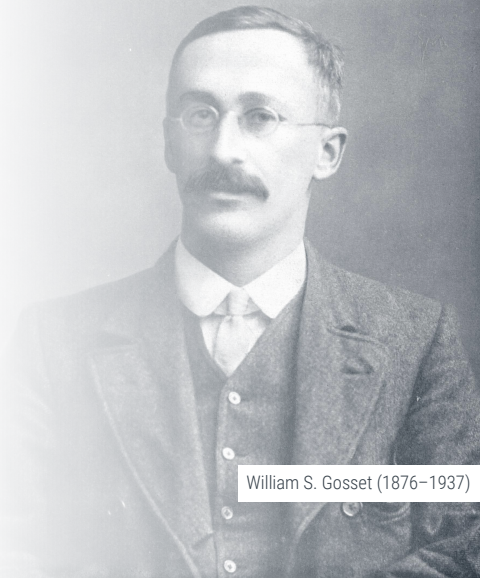▶ The conjugate prior is the *Normal-Gamma*, the predictive marginal is the *Student-t* distribution

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

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

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

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

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

William S. Gosset (1876–1937)

What if we don't know the log partition function?

▶ Consider the exponential family

$$p_w(x \mid w) = \exp\left[\phi(x)^\mathsf{T} w - \log Z(w)\right]$$

▶ for iid data:

$$p_w(x_1, x_2, \ldots, x_n \mid w) = \prod_i^n p_w(x_i \mid w) = \exp\left(\sum_i^n \phi^\mathsf{T}(x_i)w - n \log Z(w)\right)$$

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

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

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

▶ this may or may not be possible analytically, but in any case it encapsulates the data!

# Example for Maximum Likelihood Estimation

the Gaussian case

<u>Example:</u> Assume we observe samples $x_i$ drawn iid. from

$$p(x \mid \mu, \sigma^2) = \prod_{i=1}^{n} \mathcal{N}(x_i; \mu, \sigma^2) = \prod_{i} \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x_i - \mu)^2}{2\sigma^2}\right)$$

Remember that the Gaussian is an EF with

$$w := \begin{bmatrix} \frac{\mu}{\sigma^2} \\ \frac{1}{\sigma^2} \end{bmatrix} \qquad \phi(x) := \begin{bmatrix} x \\ -\frac{1}{2}x^2 \end{bmatrix} \qquad \log Z(w) := \frac{1}{2}\left(\frac{w_1^2}{w_2} - \log w_2 + \log(2\pi)\right)$$

so we find the maximum likelihood estimate by computing

$$\nabla_w \log Z(w) = \begin{bmatrix} \frac{w_1}{w_2} \\ -\frac{1}{2}\left(\frac{w_1^2}{w_2^2} + \frac{1}{w_2}\right) \end{bmatrix} = \begin{bmatrix} \mu \\ -\frac{1}{2}(\mu^2 + \sigma^2) \end{bmatrix} \stackrel{!}{=} \overline{\phi(x)} = \frac{1}{n}\sum_{i=1}^{n} \begin{bmatrix} x_i \\ -\frac{1}{2}x_i^2 \end{bmatrix}$$

hence setting $\hat{\mu} = \bar{x} =: \frac{1}{n}\sum_{i=1}^{n} x_i$ and $\hat{\sigma}^2 = \frac{1}{n}\left(\sum_{i=1}^{n} x_i^2\right) - \hat{\mu}^2$

# Why stop at maximum likelihood?

▶ thanks to auto-diff, we can do *approximate* full Bayesian inference

# Summary: Exponential Families

► reduce Bayesian inference to

  ► *modelling*: designing / computing the sufficient statistics $\phi(x)$ and partition function $Z(w)$
  ► computing the posterior: essentially evaluating the log partition function $F$ of the conjugate prior

► if $F$ is not tractable, we can still do *Laplace* approximations:

  ► find the *mode* $\hat{w}$ of the posterior, by solving the *root-finding* problem

$$\nabla_w \log p(w \mid x) = \frac{\alpha + \sum_{i=1}^{n} \phi(x_i)}{\nu + n}$$

  ► evaluate the *Hessian* $\Psi = \nabla_w \nabla_w^\intercal \log p(w \mid x)$ at $\hat{w}$
  ► approximate the posterior as $\mathcal{N}(w; \hat{w}, -\Psi^{-1})$

Laplace approximations show that Bayesian inference is not "about the prior" at all, but rather about *capturing the (local) geometry of the likelihood (around the mode)*. Uncertainty is about tracking all possible solutions at once (or at least as many as possible), not one single estimate.