

# Bayesian Machine Learning

Andrew Gordon Wilson

ORIE 6741

**Lecture 2: Bayesian Basics**

<https://people.orie.cornell.edu/andrew/orie6741>

Cornell University

August 25, 2016

# Canonical Machine Learning Problems

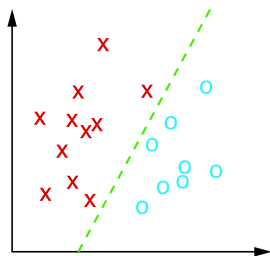
- ▶ Linear Classification
- ▶ Polynomial Regression
- ▶ Clustering with Gaussian Mixtures

# Linear Classification

► **Data:**  $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$ .

$$\mathbf{x}_i \in \mathbb{R}^D$$

$$y_i \in \{+1, -1\}.$$



► **Model:**

$$p(y_i = +1 | \boldsymbol{\theta}, \mathbf{x}) = \begin{cases} 1 & \text{if } \boldsymbol{\theta}^T \mathbf{x}_i \geq 0, \\ 0 & \text{otherwise} \end{cases} \quad (1)$$

► **Parameters:**  $\boldsymbol{\theta} \in \mathbb{R}^D$

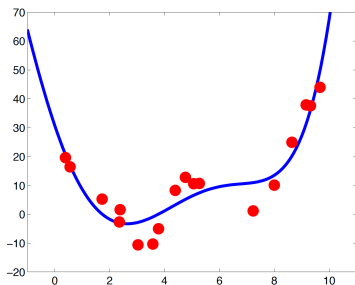
► **Goal:** To infer  $\boldsymbol{\theta}$  from data and to predict future labels  $p(y | \mathcal{D}, \mathbf{x})$ .

# Polynomial Regression

► **Data:**  $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$ .

$$\mathbf{x}_i \in \mathbb{R}$$

$$y_i \in \mathbb{R}.$$

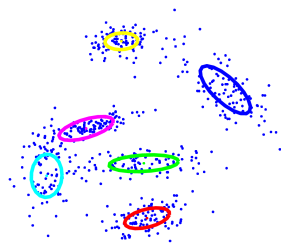


- **Model:**  $y_i = a_0 + a_1x_i + a_2x_i^2 + \dots + a_mx_i^m + \epsilon_i$ ,  
where  $\epsilon_i \sim \mathcal{N}(0, \sigma^2)$ .
- **Parameters:**  $\theta = (a_0, \dots, a_m, \sigma^2)^T$ .
- **Goal:** To infer  $\theta$  from the data and to predict future outputs  $p(y|\mathcal{D}, x, m)$ .

# Clustering with Gaussian Mixtures (Density Estimation)

► **Data:**  $\mathcal{D} = \{\mathbf{x}_i\}_{i=1}^n$

$$\mathbf{x}_i \in \mathbb{R}^D$$



► **Model:**

$$\mathbf{x}_i \sim \sum_{j=1}^m w_j p_j(\mathbf{x}_i) \quad (2)$$

$$\text{where } p_j(\mathbf{x}_i) = \mathcal{N}(\mathbf{x}_i; \boldsymbol{\mu}_j, \Sigma_j). \quad (3)$$

► **Parameters:**  $\boldsymbol{\theta} = ((\boldsymbol{\mu}_1, \Sigma_1), \dots, (\boldsymbol{\mu}_m, \Sigma_m), \mathbf{w})$

► **Goal:** To infer  $\boldsymbol{\theta}$  from the data, predict the density  $p(\mathbf{x}|\mathcal{D}, m)$ , and infer which points belong to which cluster.

# Bayesian Modelling (Theory of Everything)

*Everything follows from two simple rules:*

**Sum rule:**  $P(x) = \sum_y P(x, y)$

**Product rule:**  $P(x, y) = P(x)P(y|x)$

$$P(\theta|\mathcal{D}, m) = \frac{P(\mathcal{D}|\theta, m)P(\theta|m)}{P(\mathcal{D}|m)}$$

$P(\mathcal{D}|\theta, m)$  likelihood of parameters  $\theta$  in model  $m$   
 $P(\theta|m)$  prior probability of  $\theta$   
 $P(\theta|\mathcal{D}, m)$  posterior of  $\theta$  given data  $\mathcal{D}$

**Prediction:**

$$P(x|\mathcal{D}, m) = \int P(x|\theta, \mathcal{D}, m)P(\theta|\mathcal{D}, m)d\theta$$

**Model Comparison:**

$$P(m|\mathcal{D}) = \frac{P(\mathcal{D}|m)P(m)}{P(\mathcal{D})}$$

$$P(\mathcal{D}|m) = \int P(\mathcal{D}|\theta, m)P(\theta|m) d\theta$$

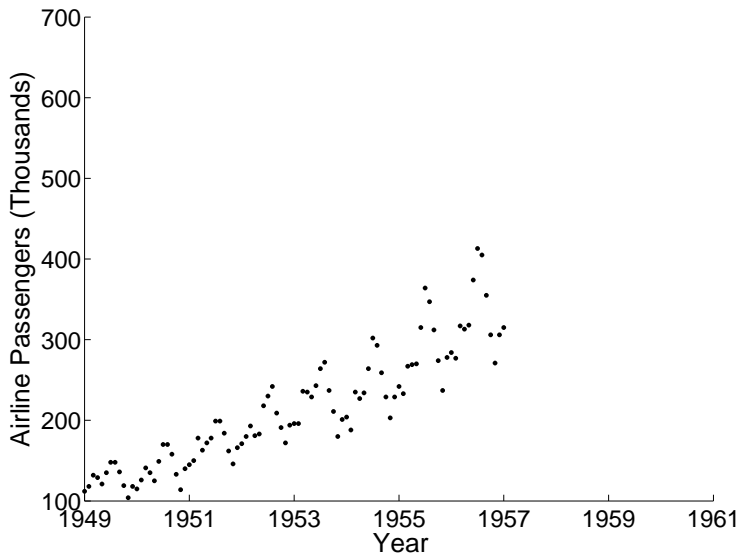
## Basic Regression Problem

- ▶ Training set of  $N$  targets (observations)  $\mathbf{y} = (y(x_1), \dots, y(x_N))^T$ .
- ▶ Observations evaluated at inputs  $X = (x_1, \dots, x_N)^T$ .
- ▶ Want to predict the value of  $y(x_*)$  at a test input  $x_*$ .

For example: Given CO<sub>2</sub> concentrations  $\mathbf{y}$  measured at times  $X$ , what will the CO<sub>2</sub> concentration be for  $x_* = 2024$ , 10 years from now?

Just knowing high school math, what might you try?

# Statistics from Scratch





## Guess the parametric form of a function that could fit the data

- ▶  $f(x, \mathbf{w}) = \mathbf{w}^T x$  [Linear function of  $\mathbf{w}$  and  $x$ ]
- ▶  $f(x, \mathbf{w}) = \mathbf{w}^T \phi(x)$  [Linear function of  $\mathbf{w}$ ] (Linear Basis Function Model)
- ▶  $f(x, \mathbf{w}) = g(\mathbf{w}^T \phi(x))$  [Non-linear in  $x$  and  $\mathbf{w}$ ] (E.g., Neural Network)

$\phi(x)$  is a vector of basis functions. For example, if  $\phi(x) = (1, x, x^2)$  and  $x \in \mathbb{R}^1$  then  $f(x, \mathbf{w}) = w_0 + w_1x + w_2x^2$  is a quadratic function.

## Choose an error measure $E(\mathbf{w})$ , minimize with respect to $\mathbf{w}$

- ▶  $E(\mathbf{w}) = \sum_{i=1}^N [f(x_i, \mathbf{w}) - y(x_i)]^2$

# Statistics from Scratch

## A probabilistic approach

We could explicitly account for noise in our model.

- ▶  $y(x) = f(x, \mathbf{w}) + \epsilon(x)$ , where  $\epsilon(x)$  is a noise function.

One commonly takes  $\epsilon(x) = \mathcal{N}(0, \sigma^2)$  for i.i.d. additive Gaussian noise, in which case

$$p(y(x)|x, \mathbf{w}, \sigma^2) = \mathcal{N}(y(x); f(x, \mathbf{w}), \sigma^2) \quad \text{Observation Model} \quad (4)$$

$$p(\mathbf{y}|x, \mathbf{w}, \sigma^2) = \prod_{i=1}^N \mathcal{N}(y(x_i); f(x_i, \mathbf{w}), \sigma^2) \quad \text{Likelihood} \quad (5)$$

- ▶ Maximize the likelihood of the data  $p(\mathbf{y}|x, \mathbf{w}, \sigma^2)$  with respect to  $\sigma^2, \mathbf{w}$ .

For a Gaussian noise model, this approach will make the same predictions as using a squared loss error function:

$$\log p(\mathbf{y}|X, \mathbf{w}, \sigma^2) \propto -\frac{1}{2\sigma^2} \sum_{i=1}^N [f(x_i, \mathbf{w}) - y(x_i)]^2 \quad (6)$$

# Statistics from Scratch

- ▶ The probabilistic approach helps us interpret the error measure in a deterministic approach, and gives us a sense of the noise level  $\sigma^2$ .
- ▶ Probabilistic methods thus provide an intuitive framework for representing uncertainty, and model development.
- ▶ Both approaches are prone to *over-fitting* for flexible  $f(x, \mathbf{w})$ : low error on the training data, high error on the test set.

## Regularization

- ▶ Use a penalized log likelihood (or error function), such as

$$\log p(\mathbf{y}|X, \mathbf{w}) \propto \underbrace{-\frac{1}{2\sigma^2} \sum_{i=1}^n (f(x_i, \mathbf{w}) - y(x_i))^2}_{\text{model fit}} \underbrace{-\lambda \mathbf{w}^T \mathbf{w}}_{\text{complexity penalty}}. \quad (7)$$

- ▶ **But how should we define complexity, and how much should we penalize complexity?**
- ▶ Can set  $\lambda$  using *cross-validation*.

## Bayes' Rule

$$p(a|b) = p(b|a)p(a)/p(b), \quad p(a|b) \propto p(b|a)p(a). \quad (8)$$

$$\text{posterior} = \frac{\text{likelihood} \times \text{prior}}{\text{marginal likelihood}}, \quad p(\mathbf{w}|\mathbf{y}, X, \sigma^2) = \frac{p(\mathbf{y}|X, \mathbf{w}, \sigma^2)p(\mathbf{w})}{p(\mathbf{y}|X, \sigma^2)}. \quad (9)$$

## Predictive Distribution

$$p(y|x_*, \mathbf{y}, X) = \int p(y|x_*, \mathbf{w})p(\mathbf{w}|\mathbf{y}, X)d\mathbf{w}. \quad (10)$$

- ▶ Think of each setting of  $\mathbf{w}$  as a different model. Eq. (10) is a *Bayesian model average*, an average of infinitely many models weighted by their posterior probabilities.
- ▶ No over-fitting, automatically calibrated complexity.
- ▶ Eq. (10) is intractable for many likelihoods  $p(\mathbf{y}|X, \mathbf{w}, \sigma^2)$  and priors  $p(\mathbf{w})$ .
- ▶ Typically more interested in the induced distribution over functions than in parameters  $\mathbf{w}$ . Can be hard to have intuitions for priors on  $p(\mathbf{w})$ .

# Parametric Regression Review

## Deterministic

$$E(\mathbf{w}) = \sum_{i=1}^N (f(x_i, \mathbf{w}) - y_i)^2. \quad (11)$$

## Maximum Likelihood

$$p(y(x)|x, \mathbf{w}) = \mathcal{N}(y(x); f(x, \mathbf{w}), \sigma_n^2), \quad (12)$$

$$p(\mathbf{y}|X, \mathbf{w}) = \prod_{i=1}^N \mathcal{N}(y(x_i); f(x_i, \mathbf{w}), \sigma_n^2). \quad (13)$$

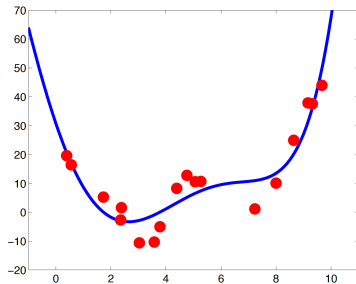
## Bayesian

$$\text{posterior} = \frac{\text{likelihood} \times \text{prior}}{\text{marginal likelihood}}, \quad p(\mathbf{w}|\mathbf{y}, X) = \frac{p(\mathbf{y}|X, \mathbf{w})p(\mathbf{w})}{p(\mathbf{y}|X)}. \quad (14)$$

$$p(y|x_*, \mathbf{y}, X) = \int p(y|x_*, \mathbf{w})p(\mathbf{w}|\mathbf{y}, X)d\mathbf{w}. \quad (15)$$

# Worked Example: Basis Regression (Chalkboard)

- ▶ We have data  $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^N$



- ▶ We use the model:

$$y = \mathbf{w}^T \phi(\mathbf{x}) + \epsilon \quad (16)$$

$$\epsilon \sim \mathcal{N}(0, \sigma^2). \quad (17)$$

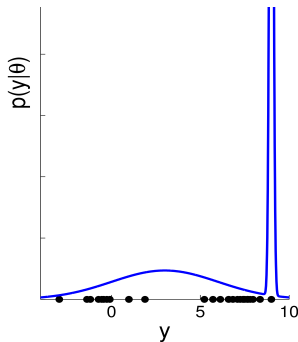
- ▶ We want to make predictions of  $y_*$  for any  $\mathbf{x}_*$ .
- ▶ We will now explore this question on the whiteboard using maximum likelihood and Bayesian approaches.
- ▶ We will consider topics such as *regularization*, *cross-validation*, *Bayesian model averaging* and *conjugate priors*.

# Rant: Regularisation = MAP $\neq$ Bayesian Inference

## Example: Density Estimation

- ▶ Observations  $y_1, \dots, y_N$  drawn from unknown density  $p(y)$ .
- ▶ Model  $p(y|\theta) = w_1\mathcal{N}(y|\mu_1, \sigma_1^2) + w_2\mathcal{N}(y|\mu_2, \sigma_2^2)$ ,  
 $\theta = \{w_1, w_2, \mu_1, \mu_2, \sigma_1, \sigma_2\}$ .
- ▶ Likelihood  $p(\mathbf{y}|\theta) = \prod_{i=1}^N p(y_i|\theta)$ .

Can learn all free parameters  $\theta$  using maximum likelihood...



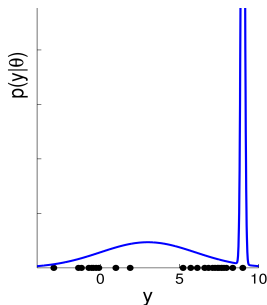
# Regularisation = MAP $\neq$ Bayesian Inference

## Regularisation or MAP

- ▶ Find

$$\operatorname{argmax}_{\theta} \log p(\theta|\mathbf{y}) \stackrel{c}{=} \underbrace{\log p(\mathbf{y}|\theta)}_{\text{model fit}} + \underbrace{\log p(\theta)}_{\text{complexity penalty}}$$

- ▶ Choose  $p(\theta)$  such that  $p(\theta) \rightarrow 0$  faster than  $p(\mathbf{y}|\theta) \rightarrow \infty$  as  $\sigma_1$  or  $\sigma_2 \rightarrow 0$ .



## Bayesian Inference

- ▶ Predictive Distribution:  $p(y_*|\mathbf{y}) = \int p(y_*|\theta)p(\theta|\mathbf{y})d\theta$ .
- ▶ Parameter Posterior:  $p(\theta|\mathbf{y}) \propto p(\mathbf{y}|\theta)p(\theta)$ .
- ▶  $p(\theta)$  need not be zero anywhere in order to make reasonable inferences. Can use a sampling scheme, with conjugate posterior updates for each separate mixture component, using an inverse Gamma prior on the variances  $\sigma_1^2, \sigma_2^2$ .



# Learning with Stochastic Gradient Descent

Chalkboard.