

Bayesian Estimation

Outline

- Review of maximum likelihood estimation
- Bayesian estimation
- Examples in machine learning

Maximum Likelihood Estimation

- A *probability model* explains how data are generated, given a parameter θ

$$\mathbf{Z}_1, \dots, \mathbf{Z}_n \stackrel{iid}{\sim} p(\mathbf{z}; \theta)$$

where $p(\mathbf{z}; \theta)$ is a pdf/pmf

- Given a realization of data $\mathbf{z}_1, \dots, \mathbf{z}_n$, the *likelihood function* is

$$L(\theta) := \prod_{i=1}^n p(\mathbf{z}_i; \theta)$$

and the *log-likelihood* is

$$\ell(\theta) := \log L(\theta) = \sum_{i=1}^n \log p(\mathbf{z}_i; \theta).$$

- A *maximum likelihood estimate* (MLE) is any

$$\hat{\theta} \in \arg \max_{\theta} \ell(\theta)$$

MLE Examples

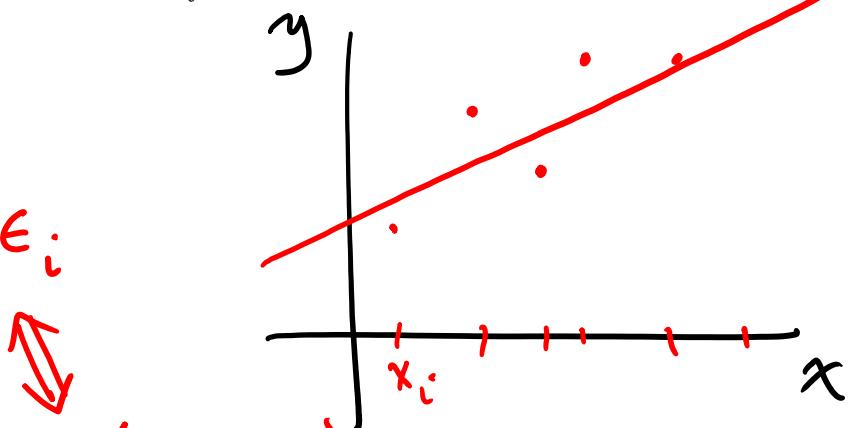
- LDA
- Naïve Bayes
- Logistic regression
- Linear regression ← *today*
- Gaussian mixture models

Maximum Likelihood Linear Regression

- Training data: $(x_1, y_1), \dots, (x_n, y_n)$. View x_i as fixed.
- No offset, $\theta = w$
- Gaussian likelihood:

$$y_i = w^T x_i + \epsilon_i$$

where $\epsilon_i \stackrel{iid}{\sim} \mathcal{N}(0, \sigma_e^2)$, σ_e^2 known.



$$y_i \sim N(w^T x_i, \sigma_e^2)$$

- Likelihood

$$L(w) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma_e^2}} \exp \left\{ -\frac{1}{2\sigma_e^2} (y_i - w^T x_i)^2 \right\}$$

- Log-likelihood

$$-\frac{1}{2\sigma_e^2} \sum_{i=1}^n (y_i - w^T x_i)^2 + \text{const.}$$

$MLE \iff OLS$

Bayesian Estimation

user-specified

- Maximum likelihood assumes that θ is nonrandom
- In Bayesian estimation, θ is viewed as **random**
- Our uncertainty about θ depends on whether we have observed data or not
- Let $\underline{Z} = (Z_1, \dots, Z_n)$ denote all observed quantities
- Let $p(\theta)$ denote the distribution of θ before observing a realization \underline{z} of \underline{Z} . **prior distribution**
- Let $p(\theta | \underline{z})$ denote the distribution of θ after observing a realization \underline{z} of \underline{Z} . **posterior distribution**
- These are related by Bayes rule:

$$p(\theta | \underline{z}) = \frac{p(\underline{z} | \theta) p(\theta)}{p(\underline{z})}$$

(\underline{z}, θ) jointly distributed

Bayesian Estimation

- A parameter estimate can be obtained from the posterior in multiple ways:

- Posterior mean: $\hat{\theta}(\underline{z}) = \mathbb{E}[\theta | \underline{z} = \underline{z}] = \int \theta p(\theta | \underline{z}) d\theta$

- Maximum a posteriori (MAP) estimation:

$$\begin{aligned}\hat{\theta}(\underline{z}) &= \arg \max_{\theta} \log p(\theta | \underline{z}) \\ &= \arg \max_{\theta} \log p(\underline{z} | \theta) + \log p(\theta)\end{aligned}$$

- In addition to parameter estimates, the main advantage of Bayesian methods over frequentist methods is they also provide natural

confidence intervals.

Bayesian Linear Regression

- Training data: $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)$. View \mathbf{x}_i as fixed.
- No offset, $\theta = \mathbf{w}$
- Gaussian likelihood:

$$y_i = \mathbf{w}^T \mathbf{x}_i + \epsilon_i, \quad \iff \quad y_i \sim N(\mathbf{w}^T \mathbf{x}_i, \sigma_e^2)$$

where $\epsilon_i \sim N(0, \sigma_e^2)$

- Prior:

$$\mathbf{w} \sim N(\mathbf{0}, \sigma_w^2 \mathbf{I}) \iff p(\mathbf{w}) = (2\pi\sigma_w^2)^{-\frac{d}{2}} \exp\left\{-\frac{\|\mathbf{w}\|^2}{2\sigma_w^2}\right\}$$

$$\iff \mathbf{w} = \begin{bmatrix} w_1 \\ \vdots \\ w_d \end{bmatrix}, \quad w_i \stackrel{iid}{\sim} N(0, \sigma_w^2)$$

- σ_e^2, σ_w^2 are hyperparameters

Bayesian Linear Regression

- Joint distribution:

$$y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}, \epsilon = \begin{bmatrix} \epsilon_1 \\ \vdots \\ \epsilon_n \end{bmatrix}, X = \begin{bmatrix} x_1 & \cdots & x_n \end{bmatrix}_{(d \times n)}$$

$$\begin{bmatrix} y \\ w \end{bmatrix} = \begin{bmatrix} I & X^T \\ 0 & I \end{bmatrix} \begin{bmatrix} \epsilon \\ w \end{bmatrix}$$

$(n+d) \times 1$

Poll

Which of the following best describes the posterior distribution of w given \hat{y} :

- (A) Gaussian
- (B) Mixture of Gaussians
- (C) Gaussian with low rank covariance
- (D) None of the above

, i.e. the
distribution

Bayesian Linear Regression

- Posterior:

$$w|y \sim N(\mu(y), \Sigma(y))$$

where

$$\hat{\omega} = \mu(y) = \left(\mathbf{X} \mathbf{X}^T + \frac{\sigma_e^2}{\sigma_w^2} \mathbf{I} \right)^{-1} \mathbf{X} y$$

$$\Sigma(y) = \left(\frac{1}{\sigma_e^2} \mathbf{X} \mathbf{X}^T + \frac{1}{\sigma_w^2} \mathbf{I} \right)^{-1}$$

and

$$\mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}, \quad \mathbf{X} = \begin{bmatrix} \mathbf{x}_1 & \cdots & \mathbf{x}_n \end{bmatrix} \in \mathbb{R}^{d \times n}$$

MAP Estimate

- Log-posterior

$$\begin{aligned}
 \log p(w|y) &= \log \frac{p(y|w)p(w)}{p(y)} \\
 &= \log p(y|w) + \log p(w) + \text{const.} \\
 &= -\frac{1}{2\sigma_e^2} \sum_{i=1}^n (y_i - w^\top x_i)^2 - \frac{1}{2\sigma_w^2} \|w\|^2 + \text{const}
 \end{aligned}$$

MAP: $\max_w p(w|y) \iff \min_w \frac{1}{2\sigma_e^2} \sum_{i=1}^n (y_i - w^\top x_i)^2 + \frac{1}{2\sigma_w^2} \|w\|^2$

- Compare with ridge regression

$$\min_w \frac{1}{n} \sum_{i=1}^n (y_i - w^\top x_i)^2 + \lambda \|w\|^2, \quad \lambda = \frac{\sigma_e^2}{n\sigma_w^2}$$

Bayesian Linear Regression

- Linear transformation of a MVG is another MVG
- In particular, if $\mathbf{Z} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, and \mathbf{A} is a matrix and \mathbf{b} a column vector of appropriate dimensions, then

$$\mathbf{AZ} + \mathbf{b} \sim \mathcal{N}(\mathbf{A}\boldsymbol{\mu} + \mathbf{b}, \mathbf{A}\boldsymbol{\Sigma}\mathbf{A}^T)$$

- If \mathbf{x} is a test point, then the predicted output

$$\hat{f}(\mathbf{x}) = \mathbf{w}^T \mathbf{x} = \mathbf{x}^T \mathbf{w}$$

- Thus, the predicted output is a random variable with distribution

$$\hat{f}(\mathbf{x}) | \mathbf{y} \sim \mathcal{N}(\mathbf{x}^T \boldsymbol{\mu}(\mathbf{y}), \mathbf{x}^T \boldsymbol{\Sigma}(\mathbf{y}) \mathbf{x})$$

$$\sim \mathcal{N}\left(\underbrace{\mathbf{x}^T (\mathbf{X}\mathbf{X}^T + \frac{\sigma_e^2}{\sigma_w^2} \mathbf{I})^{-1} \mathbf{X}_y}_{M_x(y)}, \underbrace{\mathbf{x}^T (\mathbf{X}\mathbf{X}^T + \frac{\sigma_e^2}{\sigma_w^2} \mathbf{I})^{-1} \mathbf{x}}_{\sigma_x^2(y)}\right)$$

Bayesian Linear Regression

- Using matrix identities (e.g., the matrix inversion lemma), it can be shown that

$$\hat{f}(\mathbf{x}) \mid \mathbf{y} \sim \mathcal{N}(\mu_{\mathbf{x}}(\mathbf{y}), \sigma_{\mathbf{x}}^2(\mathbf{y}))$$

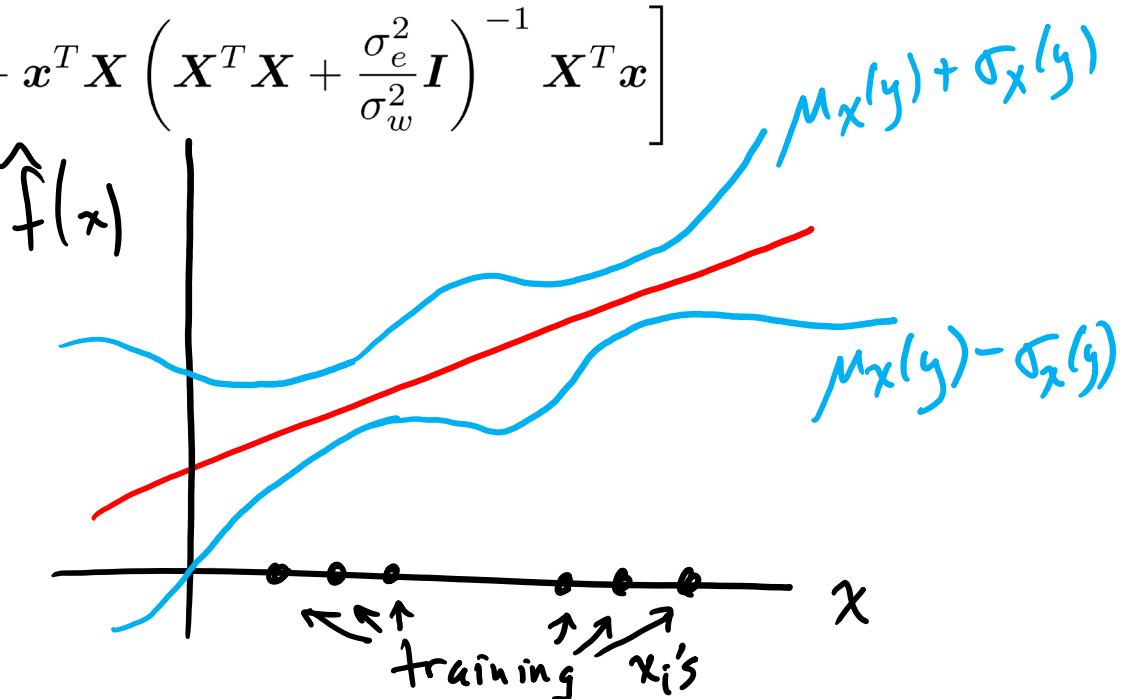
where

$$\mu_{\mathbf{x}}(\mathbf{y}) = \sigma_w^2 \mathbf{x}^T \mathbf{X} \left(\mathbf{X}^T \mathbf{X} + \frac{\sigma_e^2}{\sigma_w^2} \mathbf{I} \right)^{-1} \mathbf{y}$$

$$\sigma_{\mathbf{x}}^2(\mathbf{y}) = \sigma_w^2 \left[\mathbf{x}^T \mathbf{x} - \mathbf{x}^T \mathbf{X} \left(\mathbf{X}^T \mathbf{X} + \frac{\sigma_e^2}{\sigma_w^2} \mathbf{I} \right)^{-1} \mathbf{X}^T \mathbf{x} \right]$$

- Conclusion:

You can
kernelize
this
method

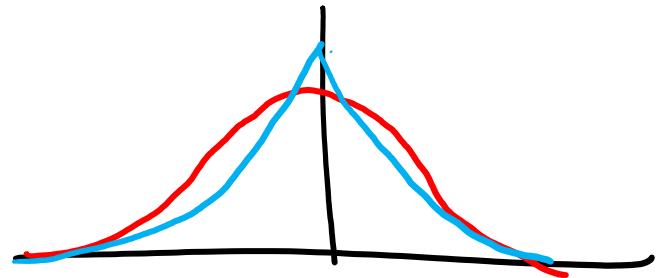


Laplacian Prior

- Training data: $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)$. View \mathbf{x}_i as fixed.
- No offset, $\theta = \mathbf{w}$
- Gaussian likelihood:

$$y_i = \mathbf{w}^T \mathbf{x}_i + \epsilon_i,$$

where $\epsilon_i \sim \mathcal{N}(0, \sigma_e^2)$



- Prior:

$$\mathbf{w} = \begin{bmatrix} w_1 \\ \vdots \\ w_d \end{bmatrix}, \quad w_i \stackrel{iid}{\sim} \text{Laplacian}(\beta)$$

$$p(w_j) = \frac{\beta}{2} \exp(-\beta |w_j|)$$

MAP Estimate

- Log-posterior

$$\begin{aligned}\log p(w|y) &= \log p(y|w) + \log p(w) + C \\ &= -\frac{1}{2\sigma_e^2} \sum_{i=1}^n (y_i - w^\top x_i)^2 - \beta \sum_{j=1}^d |w_j| + C\end{aligned}$$

$$\text{MAP} \iff \min_w \frac{1}{2\sigma_e^2} \sum_{i=1}^n (y_i - w^\top x_i)^2 + \beta \|w\|_1$$

- Compare with lasso regression

$$\min_w \frac{1}{n} \sum_{i=1}^n (y_i - w^\top x_i)^2 + \lambda \|w\|_2, \quad \lambda = \frac{\beta}{2n\sigma_e^2}$$

Poll

A penalized negative log-likelihood has the form

$$-\ell(\boldsymbol{\theta}) + \text{pen}(\boldsymbol{\theta})$$

Example: logistic regression with quadratic penalty

True or false: Every estimate obtained by minimizing a penalized negative log-likelihood can be viewed as a MAP estimate. That is, for every penalty $\text{pen}(\boldsymbol{\theta})$ there is a prior $p(\boldsymbol{\theta})$ such that $\text{pen}(\boldsymbol{\theta}) = -\log p(\boldsymbol{\theta})$.

- (A) True
- (B) False

$$1 = \int p(\boldsymbol{\theta}) d\boldsymbol{\theta} = \int e^{-\text{pen}(\boldsymbol{\theta})} d\boldsymbol{\theta}$$

no reason to equal 1

$$\Leftrightarrow p(\boldsymbol{\theta}) = e^{-\text{pen}(\boldsymbol{\theta})}$$

$$\boxed{\begin{aligned} \boldsymbol{\theta} &> 1 \\ \text{pen}(\boldsymbol{\theta}) &= \log \log \boldsymbol{\theta} \end{aligned}}$$

Multinomial Parameter Estimation

- Consider a pmf on K items given by $\theta_1, \dots, \theta_K$
- After N iid draws, let

$N_k = \# \text{ of occurrences of outcome } k$

- Then we say (N_1, \dots, N_K) have a *multinomial*

and write

$$(N_1, \dots, N_K) \sim \text{multi}(N; \theta_1, \dots, \theta_N)$$

- MLE:

$$\hat{\theta}_k = \frac{N_k}{N}$$

Multinomial Parameter Estimation

- Dirichlet prior:

$$p(\theta_1, \dots, \theta_K; \alpha_1, \dots, \alpha_K) = \frac{1}{B(\alpha)} \prod_{k=1}^K \theta_k^{\alpha_k - 1}$$

- Posterior:

$$p(\theta_1, \dots, \theta_K; \alpha_1 + N_1, \dots, \alpha_K + N_K) \quad \alpha_k = 1 \forall k \Rightarrow \text{uniform}$$

- Posterior mode (MAP estimate):

$$\hat{\theta}_i = \frac{\alpha_i}{\sum \alpha_k}$$

- Posterior mean:

$$\hat{\theta}_i = \frac{\alpha_i - 1}{\sum \alpha_k - K}$$

{

$$\alpha_i \rightarrow \alpha_i + N_i$$

Final Thoughts

- When the posterior and prior distributions belong to the same family of distributions, the prior is said to be *conjugate* to the likelihood
- Priors are typically chosen to make posterior calculations easier
- Advantages of Bayesian estimation:
 - Natural confidence intervals
 - Make use of prior knowledge if available
- Disadvantages of Bayesian estimation:
 - Inference may be computationally demanding for complex models
 - Still requires tuning hyperparameters
 - May perform poorly if prior misspecified
- Next time: Gaussian processes for regression