

Probabilistic Models

Shan-Hung Wu
shwu@cs.nthu.edu.tw

Department of Computer Science,
National Tsing Hua University, Taiwan

Machine Learning

Outline

- 1 Probabilistic Models
- 2 Maximum Likelihood Estimation
 - Linear Regression
 - Logistic Regression
- 3 Maximum A Posteriori Estimation
- 4 Bayesian Estimation**

Outline

- 1 Probabilistic Models
- 2 Maximum Likelihood Estimation
 - Linear Regression
 - Logistic Regression
- 3 Maximum A Posteriori Estimation
- 4 Bayesian Estimation**

Predictions based on Probability

- Supervised learning, we are given a training set $\mathbb{X} = \{(\mathbf{x}^{(i)}, \mathbf{y}^{(i)})\}_{i=1}^N$

Predictions based on Probability

- Supervised learning, we are given a training set $\mathbb{X} = \{(\mathbf{x}^{(i)}, \mathbf{y}^{(i)})\}_{i=1}^N$
- Model \mathbb{F} : a collection of functions parametrized by Θ

Predictions based on Probability

- Supervised learning, we are given a training set $\mathbb{X} = \{(\mathbf{x}^{(i)}, \mathbf{y}^{(i)})\}_{i=1}^N$
- Model \mathbb{F} : a collection of functions parametrized by Θ
- Goal: to train a function f such that, given a new data point \mathbf{x}' , the output value

$$\hat{\mathbf{y}} = f(\mathbf{x}'; \Theta)$$

is closest to the correct label \mathbf{y}'

Predictions based on Probability

- Supervised learning, we are given a training set $\mathbb{X} = \{(\mathbf{x}^{(i)}, \mathbf{y}^{(i)})\}_{i=1}^N$
- Model \mathbb{F} : a collection of functions parametrized by Θ
- Goal: to train a function f such that, given a new data point \mathbf{x}' , the output value

$$\hat{\mathbf{y}} = f(\mathbf{x}'; \Theta)$$

is closest to the correct label \mathbf{y}'

- Examples in \mathbb{X} are usually assumed to be i.i.d. sampled from random variables (\mathbf{x}, \mathbf{y}) following some data generating distribution $P(\mathbf{x}, \mathbf{y})$

Predictions based on Probability

- Supervised learning, we are given a training set $\mathbb{X} = \{(\mathbf{x}^{(i)}, \mathbf{y}^{(i)})\}_{i=1}^N$
- Model \mathbb{F} : a collection of functions parametrized by Θ
- Goal: to train a function f such that, given a new data point \mathbf{x}' , the output value

$$\hat{\mathbf{y}} = f(\mathbf{x}'; \Theta)$$

is closest to the correct label \mathbf{y}'

- Examples in \mathbb{X} are usually assumed to be i.i.d. sampled from random variables (\mathbf{x}, \mathbf{y}) following some data generating distribution $P(\mathbf{x}, \mathbf{y})$
- In probabilistic models, f is replaced by $P(\mathbf{y} = \mathbf{y}' | \mathbf{x} = \mathbf{x}')$ and a prediction is made by:

$$\hat{\mathbf{y}} = \arg \max_{\mathbf{y}} P(\mathbf{y} = \mathbf{y}' | \mathbf{x} = \mathbf{x}'; \Theta)$$

Predictions based on Probability

- Supervised learning, we are given a training set $\mathbb{X} = \{(\mathbf{x}^{(i)}, \mathbf{y}^{(i)})\}_{i=1}^N$
- Model \mathbb{F} : a collection of functions parametrized by Θ
- Goal: to train a function f such that, given a new data point \mathbf{x}' , the output value

$$\hat{\mathbf{y}} = f(\mathbf{x}'; \Theta)$$

is closest to the correct label \mathbf{y}'

- Examples in \mathbb{X} are usually assumed to be i.i.d. sampled from random variables (\mathbf{x}, \mathbf{y}) following some data generating distribution $P(\mathbf{x}, \mathbf{y})$
- In probabilistic models, f is replaced by $P(\mathbf{y} = \mathbf{y}' | \mathbf{x} = \mathbf{x}')$ and a prediction is made by:

$$\hat{\mathbf{y}} = \arg \max_{\mathbf{y}} P(\mathbf{y} = \mathbf{y}' | \mathbf{x} = \mathbf{x}'; \Theta)$$

- How to find Θ ?

Function (Θ) as Point Estimate

- Regard $\Theta(f)$ as an estimate of the “true” $\Theta^*(f^*)$
 - Mapped from the training set \mathbb{X}

Function (Θ) as Point Estimate

- Regard $\Theta(f)$ as an estimate of the “true” $\Theta^*(f^*)$
 - Mapped from the training set \mathbb{X}
- **Maximum a posteriori (MAP) estimation:**

$$\arg \max_{\Theta} P(\Theta | \mathbb{X}) = \arg \max_{\Theta} P(\mathbb{X} | \Theta)P(\Theta)$$

- By Bayes' rule ($P(\mathbb{X})$ is irrelevant)

Function (Θ) as Point Estimate

- Regard $\Theta(f)$ as an estimate of the “true” $\Theta^*(f^*)$
 - Mapped from the training set \mathbb{X}
- **Maximum a posteriori (MAP) estimation:**

$$\arg \max_{\Theta} P(\Theta | \mathbb{X}) = \arg \max_{\Theta} P(\mathbb{X} | \Theta)P(\Theta)$$

- By Bayes' rule ($P(\mathbb{X})$ is irrelevant)
- Solves Θ first, then uses it as a constant in $P(\mathbf{y}|\mathbf{x}; \Theta)$ to get $\hat{\mathbf{y}}$

Function (Θ) as Point Estimate

- Regard $\Theta(f)$ as an estimate of the “true” $\Theta^*(f^*)$
 - Mapped from the training set \mathbb{X}
- **Maximum a posteriori (MAP) estimation:**

$$\arg \max_{\Theta} P(\Theta | \mathbb{X}) = \arg \max_{\Theta} P(\mathbb{X} | \Theta)P(\Theta)$$

- By Bayes' rule ($P(\mathbb{X})$ is irrelevant)
 - Solves Θ first, then uses it as a constant in $P(\mathbf{y} | \mathbf{x}; \Theta)$ to get $\hat{\mathbf{y}}$
- **Maximum likelihood (ML) estimation:**

$$\arg \max_{\Theta} P(\mathbb{X} | \Theta)$$

- Assumes uniform $P(\Theta)$ and does not prefer particular Θ

Outline

- 1 Probabilistic Models
- 2 Maximum Likelihood Estimation**
 - Linear Regression
 - Logistic Regression
- 3 Maximum A Posteriori Estimation
- 4 Bayesian Estimation**

Outline

- 1 Probabilistic Models
- 2 Maximum Likelihood Estimation**
 - Linear Regression
 - Logistic Regression
- 3 Maximum A Posteriori Estimation
- 4 Bayesian Estimation**

Probability Interpretation

- Assumption: $y = f^*(\mathbf{x}) + \varepsilon$, $\varepsilon \sim \mathcal{N}(0, \beta^{-1})$

Probability Interpretation

- Assumption: $y = f^*(\mathbf{x}) + \varepsilon$, $\varepsilon \sim \mathcal{N}(0, \beta^{-1})$
- The unknown deterministic function is defined as

$$f^*(\mathbf{x}; \mathbf{w}^*) = \mathbf{w}^{*\top} \mathbf{x}$$

- All variables are z -normalized, so no bias term (b)

Probability Interpretation

- Assumption: $y = f^*(\mathbf{x}) + \varepsilon$, $\varepsilon \sim \mathcal{N}(0, \beta^{-1})$
- The unknown deterministic function is defined as

$$f^*(\mathbf{x}; \mathbf{w}^*) = \mathbf{w}^{*\top} \mathbf{x}$$

- All variables are z -normalized, so no bias term (b)
- We have $(y | \mathbf{x}) \sim \mathcal{N}(\mathbf{w}^{*\top} \mathbf{x}, \beta^{-1})$

Probability Interpretation

- Assumption: $y = f^*(\mathbf{x}) + \varepsilon$, $\varepsilon \sim \mathcal{N}(0, \beta^{-1})$
- The unknown deterministic function is defined as

$$f^*(\mathbf{x}; \mathbf{w}^*) = \mathbf{w}^{*\top} \mathbf{x}$$

- All variables are z -normalized, so no bias term (b)
- We have $(y | \mathbf{x}) \sim \mathcal{N}(\mathbf{w}^{*\top} \mathbf{x}, \beta^{-1})$
- So, our goal is to find \mathbf{w} as close to \mathbf{w}^* as possible such that:

$$\hat{y} = \arg \max_y P(y | \mathbf{x} = \mathbf{x}; \mathbf{w}) = \mathbf{w}^\top \mathbf{x}$$

- Note that \hat{y} is irrelevant to β , so we don't need to solve β

Probability Interpretation

- Assumption: $y = f^*(\mathbf{x}) + \varepsilon$, $\varepsilon \sim \mathcal{N}(0, \beta^{-1})$
- The unknown deterministic function is defined as

$$f^*(\mathbf{x}; \mathbf{w}^*) = \mathbf{w}^{*\top} \mathbf{x}$$

- All variables are z -normalized, so no bias term (b)
- We have $(y | \mathbf{x}) \sim \mathcal{N}(\mathbf{w}^{*\top} \mathbf{x}, \beta^{-1})$
- So, our goal is to find \mathbf{w} as close to \mathbf{w}^* as possible such that:

$$\hat{y} = \arg \max_y P(y | \mathbf{x} = \mathbf{x}; \mathbf{w}) = \mathbf{w}^\top \mathbf{x}$$

- Note that \hat{y} is irrelevant to β , so we don't need to solve β
- ML estimation:

$$\arg \max_{\mathbf{w}} P(\mathbb{X} | \mathbf{w})$$

ML Estimation I

- Problem:

$$\arg \max_{\mathbf{w}} \mathbb{P}(\mathbb{X} | \mathbf{w})$$

- Since we assume i.i.d. samples, we have

$$\mathbb{P}(\mathbb{X} | \mathbf{w}) = \prod_{i=1}^N \mathbb{P}(\mathbf{x}^{(i)}, y^{(i)} | \mathbf{w}) = \prod_{i=1}^N \mathbb{P}(y^{(i)} | \mathbf{x}^{(i)}, \mathbf{w}) \mathbb{P}(\mathbf{x}^{(i)} | \mathbf{w})$$

ML Estimation I

- Problem:

$$\arg \max_{\mathbf{w}} \mathbb{P}(\mathbb{X} | \mathbf{w})$$

- Since we assume i.i.d. samples, we have

$$\begin{aligned} \mathbb{P}(\mathbb{X} | \mathbf{w}) &= \prod_{i=1}^N \mathbb{P}(\mathbf{x}^{(i)}, y^{(i)} | \mathbf{w}) = \prod_{i=1}^N \mathbb{P}(y^{(i)} | \mathbf{x}^{(i)}, \mathbf{w}) \mathbb{P}(\mathbf{x}^{(i)} | \mathbf{w}) \\ &= \prod_{i=1}^N \mathbb{P}(y^{(i)} | \mathbf{x}^{(i)}, \mathbf{w}) \mathbb{P}(\mathbf{x}^{(i)}) = \prod_i \mathcal{N}(y^{(i)}; \mathbf{w}^\top \mathbf{x}^{(i)}, \sigma^2) \mathbb{P}(\mathbf{x}^{(i)}) \end{aligned}$$

ML Estimation I

- Problem:

$$\arg \max_{\mathbf{w}} \mathbb{P}(\mathbb{X} | \mathbf{w})$$

- Since we assume i.i.d. samples, we have

$$\begin{aligned} \mathbb{P}(\mathbb{X} | \mathbf{w}) &= \prod_{i=1}^N \mathbb{P}(\mathbf{x}^{(i)}, y^{(i)} | \mathbf{w}) = \prod_{i=1}^N \mathbb{P}(y^{(i)} | \mathbf{x}^{(i)}, \mathbf{w}) \mathbb{P}(\mathbf{x}^{(i)} | \mathbf{w}) \\ &= \prod_{i=1}^N \mathbb{P}(y^{(i)} | \mathbf{x}^{(i)}, \mathbf{w}) \mathbb{P}(\mathbf{x}^{(i)}) = \prod_i \mathcal{N}(y^{(i)}; \mathbf{w}^\top \mathbf{x}^{(i)}, \sigma^2) \mathbb{P}(\mathbf{x}^{(i)}) \\ &= \prod_i \sqrt{\frac{\beta}{2\pi}} \exp\left(-\frac{\beta}{2}(y^{(i)} - \mathbf{w}^\top \mathbf{x}^{(i)})^2\right) \mathbb{P}(\mathbf{x}^{(i)}) \end{aligned}$$

ML Estimation I

- Problem:

$$\arg \max_{\mathbf{w}} \mathbb{P}(\mathbb{X} | \mathbf{w})$$

- Since we assume i.i.d. samples, we have

$$\begin{aligned} \mathbb{P}(\mathbb{X} | \mathbf{w}) &= \prod_{i=1}^N \mathbb{P}(\mathbf{x}^{(i)}, y^{(i)} | \mathbf{w}) = \prod_{i=1}^N \mathbb{P}(y^{(i)} | \mathbf{x}^{(i)}, \mathbf{w}) \mathbb{P}(\mathbf{x}^{(i)} | \mathbf{w}) \\ &= \prod_{i=1}^N \mathbb{P}(y^{(i)} | \mathbf{x}^{(i)}, \mathbf{w}) \mathbb{P}(\mathbf{x}^{(i)}) = \prod_i \mathcal{N}(y^{(i)}; \mathbf{w}^\top \mathbf{x}^{(i)}, \sigma^2) \mathbb{P}(\mathbf{x}^{(i)}) \\ &= \prod_i \sqrt{\frac{\beta}{2\pi}} \exp\left(-\frac{\beta}{2}(y^{(i)} - \mathbf{w}^\top \mathbf{x}^{(i)})^2\right) \mathbb{P}(\mathbf{x}^{(i)}) \end{aligned}$$

- To make the problem tractable, we prefer “sums” over “products”

ML Estimation I

- Problem:

$$\arg \max_{\mathbf{w}} \mathbb{P}(\mathbb{X} | \mathbf{w})$$

- Since we assume i.i.d. samples, we have

$$\begin{aligned} \mathbb{P}(\mathbb{X} | \mathbf{w}) &= \prod_{i=1}^N \mathbb{P}(\mathbf{x}^{(i)}, y^{(i)} | \mathbf{w}) = \prod_{i=1}^N \mathbb{P}(y^{(i)} | \mathbf{x}^{(i)}, \mathbf{w}) \mathbb{P}(\mathbf{x}^{(i)} | \mathbf{w}) \\ &= \prod_{i=1}^N \mathbb{P}(y^{(i)} | \mathbf{x}^{(i)}, \mathbf{w}) \mathbb{P}(\mathbf{x}^{(i)}) = \prod_i \mathcal{N}(y^{(i)}; \mathbf{w}^\top \mathbf{x}^{(i)}, \sigma^2) \mathbb{P}(\mathbf{x}^{(i)}) \\ &= \prod_i \sqrt{\frac{\beta}{2\pi}} \exp\left(-\frac{\beta}{2}(y^{(i)} - \mathbf{w}^\top \mathbf{x}^{(i)})^2\right) \mathbb{P}(\mathbf{x}^{(i)}) \end{aligned}$$

- To make the problem tractable, we prefer “sums” over “products”
- We can instead maximize the **log likelihood**

$$\arg \max_{\mathbf{w}} \log \mathbb{P}(\mathbb{X} | \mathbf{w})$$

- The optimal point does not change since log is monotone increasing

ML Estimation I

- Problem:

$$\arg \max_{\mathbf{w}} \mathbb{P}(\mathbb{X} | \mathbf{w})$$

- Since we assume i.i.d. samples, we have

$$\begin{aligned} \mathbb{P}(\mathbb{X} | \mathbf{w}) &= \prod_{i=1}^N \mathbb{P}(\mathbf{x}^{(i)}, y^{(i)} | \mathbf{w}) = \prod_{i=1}^N \mathbb{P}(y^{(i)} | \mathbf{x}^{(i)}, \mathbf{w}) \mathbb{P}(\mathbf{x}^{(i)} | \mathbf{w}) \\ &= \prod_{i=1}^N \mathbb{P}(y^{(i)} | \mathbf{x}^{(i)}, \mathbf{w}) \mathbb{P}(\mathbf{x}^{(i)}) = \prod_i \mathcal{N}(y^{(i)}; \mathbf{w}^\top \mathbf{x}^{(i)}, \sigma^2) \mathbb{P}(\mathbf{x}^{(i)}) \\ &= \prod_i \sqrt{\frac{\beta}{2\pi}} \exp\left(-\frac{\beta}{2}(y^{(i)} - \mathbf{w}^\top \mathbf{x}^{(i)})^2\right) \mathbb{P}(\mathbf{x}^{(i)}) \end{aligned}$$

- To make the problem tractable, we prefer “sums” over “products”
- We can instead maximize the **log likelihood**

$$\begin{aligned} \arg \max_{\mathbf{w}} \log \mathbb{P}(\mathbb{X} | \mathbf{w}) \\ = \arg \max_{\mathbf{w}} \log \left[\prod_i \sqrt{\frac{\beta}{2\pi}} \exp\left(-\frac{\beta}{2}(y^{(i)} - \mathbf{w}^\top \mathbf{x}^{(i)})^2\right) \mathbb{P}(\mathbf{x}^{(i)}) \right] \end{aligned}$$

- The optimal point does not change since log is monotone increasing

ML Estimation I

- Problem:

$$\arg \max_{\mathbf{w}} \mathbb{P}(\mathbb{X} | \mathbf{w})$$

- Since we assume i.i.d. samples, we have

$$\begin{aligned} \mathbb{P}(\mathbb{X} | \mathbf{w}) &= \prod_{i=1}^N \mathbb{P}(\mathbf{x}^{(i)}, y^{(i)} | \mathbf{w}) = \prod_{i=1}^N \mathbb{P}(y^{(i)} | \mathbf{x}^{(i)}, \mathbf{w}) \mathbb{P}(\mathbf{x}^{(i)} | \mathbf{w}) \\ &= \prod_{i=1}^N \mathbb{P}(y^{(i)} | \mathbf{x}^{(i)}, \mathbf{w}) \mathbb{P}(\mathbf{x}^{(i)}) = \prod_i \mathcal{N}(y^{(i)}; \mathbf{w}^\top \mathbf{x}^{(i)}, \sigma^2) \mathbb{P}(\mathbf{x}^{(i)}) \\ &= \prod_i \sqrt{\frac{\beta}{2\pi}} \exp\left(-\frac{\beta}{2}(y^{(i)} - \mathbf{w}^\top \mathbf{x}^{(i)})^2\right) \mathbb{P}(\mathbf{x}^{(i)}) \end{aligned}$$

- To make the problem tractable, we prefer “sums” over “products”
- We can instead maximize the **log likelihood**

$$\begin{aligned} &\arg \max_{\mathbf{w}} \log \mathbb{P}(\mathbb{X} | \mathbf{w}) \\ &= \arg \max_{\mathbf{w}} \log \left[\prod_i \sqrt{\frac{\beta}{2\pi}} \exp\left(-\frac{\beta}{2}(y^{(i)} - \mathbf{w}^\top \mathbf{x}^{(i)})^2\right) \mathbb{P}(\mathbf{x}^{(i)}) \right] \\ &= \arg \max_{\mathbf{w}} N \sqrt{\frac{\beta}{2\pi}} - \frac{\beta}{2} \sum_i (y^{(i)} - \mathbf{w}^\top \mathbf{x}^{(i)})^2 + \sum_i \log \mathbb{P}(\mathbf{x}^{(i)}) \end{aligned}$$

- The optimal point does not change since log is monotone increasing

ML Estimation II

$$\arg \max_{\mathbf{w}} N \sqrt{\frac{\beta}{2\pi}} - \frac{\beta}{2} \sum_i (y^{(i)} - \mathbf{w}^\top \mathbf{x}^{(i)})^2 + \sum_i P(\mathbf{x}^{(i)})$$

- Ignoring terms irrelevant to \mathbf{w} , we have

$$\arg \min_{\mathbf{w}} \sum_i (y^{(i)} - \mathbf{w}^\top \mathbf{x}^{(i)})^2$$

ML Estimation II

$$\arg \max_{\mathbf{w}} N \sqrt{\frac{\beta}{2\pi}} - \frac{\beta}{2} \sum_i (y^{(i)} - \mathbf{w}^\top \mathbf{x}^{(i)})^2 + \sum_i P(\mathbf{x}^{(i)})$$

- Ignoring terms irrelevant to \mathbf{w} , we have

$$\arg \min_{\mathbf{w}} \sum_i (y^{(i)} - \mathbf{w}^\top \mathbf{x}^{(i)})^2$$

- In other words, we seek for \mathbf{w} by *minimizing the SSE* (sum of square errors), as we have done before
 - By, e.g., the stochastic gradient descent algorithm

ML Estimation II

$$\arg \max_{\mathbf{w}} N \sqrt{\frac{\beta}{2\pi}} - \frac{\beta}{2} \sum_i (y^{(i)} - \mathbf{w}^\top \mathbf{x}^{(i)})^2 + \sum_i P(\mathbf{x}^{(i)})$$

- Ignoring terms irrelevant to \mathbf{w} , we have

$$\arg \min_{\mathbf{w}} \sum_i (y^{(i)} - \mathbf{w}^\top \mathbf{x}^{(i)})^2$$

- In other words, we seek for \mathbf{w} by *minimizing the SSE* (sum of square errors), as we have done before
 - By, e.g., the stochastic gradient descent algorithm
- This new perspective explains our ad hoc choice of SSE for empirical risk minimization
 - Checking assumptions helps understand when model works the best

ML Estimation II

$$\arg \max_{\mathbf{w}} N \sqrt{\frac{\beta}{2\pi}} - \frac{\beta}{2} \sum_i (y^{(i)} - \mathbf{w}^\top \mathbf{x}^{(i)})^2 + \sum_i P(\mathbf{x}^{(i)})$$

- Ignoring terms irrelevant to \mathbf{w} , we have

$$\arg \min_{\mathbf{w}} \sum_i (y^{(i)} - \mathbf{w}^\top \mathbf{x}^{(i)})^2$$

- In other words, we seek for \mathbf{w} by *minimizing the SSE* (sum of square errors), as we have done before
 - By, e.g., the stochastic gradient descent algorithm
- This new perspective explains our ad hoc choice of SSE for empirical risk minimization
 - Checking assumptions helps understand when model works the best
- Also motivates new models

ML Estimation II

$$\arg \max_{\mathbf{w}} N \sqrt{\frac{\beta}{2\pi}} - \frac{\beta}{2} \sum_i (y^{(i)} - \mathbf{w}^\top \mathbf{x}^{(i)})^2 + \sum_i P(\mathbf{x}^{(i)})$$

- Ignoring terms irrelevant to \mathbf{w} , we have

$$\arg \min_{\mathbf{w}} \sum_i (y^{(i)} - \mathbf{w}^\top \mathbf{x}^{(i)})^2$$

- In other words, we seek for \mathbf{w} by *minimizing the SSE* (sum of square errors), as we have done before
 - By, e.g., the stochastic gradient descent algorithm
- This new perspective explains our ad hoc choice of SSE for empirical risk minimization
 - Checking assumptions helps understand when model works the best
- Also motivates new models. Probabilistic model for classification?

Outline

- 1 Probabilistic Models
- 2 Maximum Likelihood Estimation**
 - Linear Regression
 - Logistic Regression
- 3 Maximum A Posteriori Estimation
- 4 Bayesian Estimation**

Probabilistic Models for Binary Classification

- Probabilistic models:

$$\hat{y} = \arg \max_y P(y | \mathbf{x}; \Theta)$$

Probabilistic Models for Binary Classification

- Probabilistic models:

$$\hat{y} = \arg \max_y P(y | \mathbf{x}; \Theta)$$

- In regression, we assume $(y | \mathbf{x}) \sim \mathcal{N}$ (based on $y = f^*(\mathbf{x}) + \varepsilon$)

Probabilistic Models for Binary Classification

- Probabilistic models:

$$\hat{y} = \arg \max_y P(y | \mathbf{x}; \Theta)$$

- In regression, we assume $(y | \mathbf{x}) \sim \mathcal{N}$ (based on $y = f^*(\mathbf{x}) + \varepsilon$)
- However, Gaussian distribution is **not** applicable to binary classification
 - The values of y should concentrate in either 1 or -1

Probabilistic Models for Binary Classification

- Probabilistic models:

$$\hat{y} = \arg \max_y P(y | \mathbf{x}; \Theta)$$

- In regression, we assume $(y | \mathbf{x}) \sim \mathcal{N}$ (based on $y = f^*(\mathbf{x}) + \varepsilon$)
- However, Gaussian distribution is **not** applicable to binary classification
 - The values of y should concentrate in either 1 or -1
- Which distribution to assume?

Probabilistic Models for Binary Classification

- Probabilistic models:

$$\hat{y} = \arg \max_y P(y | \mathbf{x}; \Theta)$$

- In regression, we assume $(y | \mathbf{x}) \sim \mathcal{N}$ (based on $y = f^*(\mathbf{x}) + \varepsilon$)
- However, Gaussian distribution is **not** applicable to binary classification
 - The values of y should concentrate in either 1 or -1
- Which distribution to assume?
- Coin flipping: $(y | \mathbf{x}) \sim \text{Bernoulli}(\rho)$, where

$$P(y | \mathbf{x}; \rho) = \rho^{y'} (1 - \rho)^{(1-y')}, \text{ where } y' = \frac{y+1}{2}$$

Probabilistic Models for Binary Classification

- Probabilistic models:

$$\hat{y} = \arg \max_y P(y | \mathbf{x}; \Theta)$$

- In regression, we assume $(y | \mathbf{x}) \sim \mathcal{N}$ (based on $y = f^*(\mathbf{x}) + \varepsilon$)
- However, Gaussian distribution is **not** applicable to binary classification
 - The values of y should concentrate in either 1 or -1
- Which distribution to assume?
- Coin flipping: $(y | \mathbf{x}) \sim \text{Bernoulli}(\rho)$, where

$$P(y | \mathbf{x}; \rho) = \rho^{y'} (1 - \rho)^{(1-y')}, \text{ where } y' = \frac{y+1}{2}$$

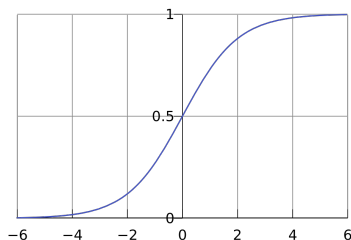
- How to relate \mathbf{x} to ρ ?

Logistic Function

- Recall that the *logistic function*

$$\sigma(z) = \frac{\exp(z)}{\exp(z) + 1} = \frac{1}{1 + \exp(-z)}$$

is commonly used as a parametrizing function of the Bernoulli distribution



Logistic Function

- Recall that the *logistic function*

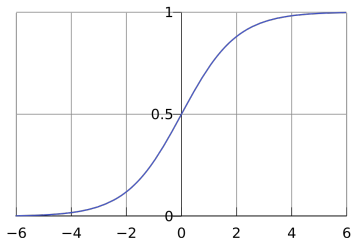
$$\sigma(z) = \frac{\exp(z)}{\exp(z) + 1} = \frac{1}{1 + \exp(-z)}$$

is commonly used as a parametrizing function of the Bernoulli distribution

- We have

$$P(y|\mathbf{x};z) = \sigma(z)^{y'} (1 - \sigma(z))^{(1-y')}$$

- The larger z , the higher chance we get a “positive flip”



Logistic Function

- Recall that the *logistic function*

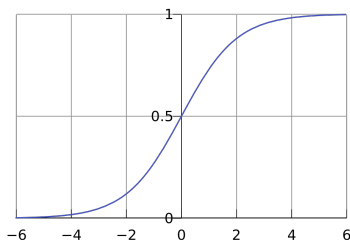
$$\sigma(z) = \frac{\exp(z)}{\exp(z) + 1} = \frac{1}{1 + \exp(-z)}$$

is commonly used as a parametrizing function of the Bernoulli distribution

- We have

$$P(y|\mathbf{x};z) = \sigma(z)^{y'} (1 - \sigma(z))^{(1-y')}$$

- The larger z , the higher chance we get a “positive flip”
- How to relate \mathbf{x} to z ?



Logistic Regression

- In *logistic regression*, we let

$$z = \mathbf{w}^\top \mathbf{x}$$

- Basically, z is the projection of \mathbf{x} along the direction \mathbf{w}

Logistic Regression

- In *logistic regression*, we let

$$z = \mathbf{w}^\top \mathbf{x}$$

- Basically, z is the projection of \mathbf{x} along the direction \mathbf{w}
- We have

$$P(y|\mathbf{x}; \mathbf{w}) = \sigma(\mathbf{w}^\top \mathbf{x})^{y'} [1 - \sigma(\mathbf{w}^\top \mathbf{x})]^{(1-y')}$$

- Prediction:

$$\hat{y} = \arg \max_y P(y|\mathbf{x}; \mathbf{w})$$

Logistic Regression

- In *logistic regression*, we let

$$z = \mathbf{w}^\top \mathbf{x}$$

- Basically, z is the projection of \mathbf{x} along the direction \mathbf{w}
- We have

$$P(y|\mathbf{x}; \mathbf{w}) = \sigma(\mathbf{w}^\top \mathbf{x})^{y'} [1 - \sigma(\mathbf{w}^\top \mathbf{x})]^{(1-y')}$$

- Prediction:

$$\hat{y} = \arg \max_y P(y|\mathbf{x}; \mathbf{w}) = \text{sign}(\mathbf{w}^\top \mathbf{x})$$

Logistic Regression

- In *logistic regression*, we let

$$z = \mathbf{w}^\top \mathbf{x}$$

- Basically, z is the projection of \mathbf{x} along the direction \mathbf{w}
- We have

$$P(y|\mathbf{x}; \mathbf{w}) = \sigma(\mathbf{w}^\top \mathbf{x})^{y'} [1 - \sigma(\mathbf{w}^\top \mathbf{x})]^{(1-y')}$$

- Prediction:

$$\hat{y} = \arg \max_y P(y|\mathbf{x}; \mathbf{w}) = \text{sign}(\mathbf{w}^\top \mathbf{x})$$

- How to learn \mathbf{w} from \mathbb{X} ?

Logistic Regression

- In *logistic regression*, we let

$$z = \mathbf{w}^\top \mathbf{x}$$

- Basically, z is the projection of \mathbf{x} along the direction \mathbf{w}
- We have

$$P(y|\mathbf{x}; \mathbf{w}) = \sigma(\mathbf{w}^\top \mathbf{x})^{y'} [1 - \sigma(\mathbf{w}^\top \mathbf{x})]^{(1-y')}$$

- Prediction:

$$\hat{y} = \arg \max_y P(y|\mathbf{x}; \mathbf{w}) = \text{sign}(\mathbf{w}^\top \mathbf{x})$$

- How to learn \mathbf{w} from \mathbb{X} ?
- ML estimation:

$$\arg \max_{\mathbf{w}} P(\mathbb{X} | \mathbf{w})$$

ML Estimation

- Log-likelihood:

$$\begin{aligned}\log P(\mathbb{X} | \mathbf{w}) &= \log \prod_{i=1}^N P(\mathbf{x}^{(i)}, y^{(i)} | \mathbf{w}) \\ &= \log \prod_i P(y^{(i)} | \mathbf{x}^{(i)}, \mathbf{w}) P(\mathbf{x}^{(i)} | \mathbf{w})\end{aligned}$$

ML Estimation

- Log-likelihood:

$$\begin{aligned}\log P(\mathbb{X} | \mathbf{w}) &= \log \prod_{i=1}^N P(\mathbf{x}^{(i)}, y^{(i)} | \mathbf{w}) \\ &= \log \prod_i P(y^{(i)} | \mathbf{x}^{(i)}, \mathbf{w}) P(\mathbf{x}^{(i)} | \mathbf{w}) \\ &\propto \log \prod_i \sigma(\mathbf{w}^\top \mathbf{x}^{(i)})^{y^{(i)}} [1 - \sigma(\mathbf{w}^\top \mathbf{x}^{(i)})]^{(1-y^{(i)})}\end{aligned}$$

ML Estimation

- Log-likelihood:

$$\begin{aligned}\log P(\mathbb{X} | \mathbf{w}) &= \log \prod_{i=1}^N P(\mathbf{x}^{(i)}, y^{(i)} | \mathbf{w}) \\ &= \log \prod_i P(y^{(i)} | \mathbf{x}^{(i)}, \mathbf{w}) P(\mathbf{x}^{(i)} | \mathbf{w}) \\ &\propto \log \prod_i \sigma(\mathbf{w}^\top \mathbf{x}^{(i)})^{y^{(i)}} [1 - \sigma(\mathbf{w}^\top \mathbf{x}^{(i)})]^{(1-y^{(i)})} \\ &= \sum_i y^{(i)} \mathbf{w}^\top \mathbf{x}^{(i)} - \log(1 + e^{\mathbf{w}^\top \mathbf{x}^{(i)}}) \text{ [Homework]}\end{aligned}$$

- Unlike in linear regression, we cannot solve \mathbf{w} analytically in a closed form via

$$\nabla_{\mathbf{w}} \log P(\mathbb{X} | \mathbf{w}) = \sum_{t=1}^N [y^{(t)} - \sigma(\mathbf{w}^\top \mathbf{x}^{(t)})] \mathbf{x}^{(t)} = \mathbf{0}$$

ML Estimation

- Log-likelihood:

$$\begin{aligned}\log P(\mathbb{X} | \mathbf{w}) &= \log \prod_{i=1}^N P(\mathbf{x}^{(i)}, y^{(i)} | \mathbf{w}) \\ &= \log \prod_i P(y^{(i)} | \mathbf{x}^{(i)}, \mathbf{w}) P(\mathbf{x}^{(i)} | \mathbf{w}) \\ &\propto \log \prod_i \sigma(\mathbf{w}^\top \mathbf{x}^{(i)})^{y^{(i)}} [1 - \sigma(\mathbf{w}^\top \mathbf{x}^{(i)})]^{(1-y^{(i)})} \\ &= \sum_i y^{(i)} \mathbf{w}^\top \mathbf{x}^{(i)} - \log(1 + e^{\mathbf{w}^\top \mathbf{x}^{(i)}}) \text{ [Homework]}\end{aligned}$$

- Unlike in linear regression, we cannot solve \mathbf{w} analytically in a closed form via

$$\nabla_{\mathbf{w}} \log P(\mathbb{X} | \mathbf{w}) = \sum_{t=1}^N [y^{(t)} - \sigma(\mathbf{w}^\top \mathbf{x}^{(t)})] \mathbf{x}^{(t)} = \mathbf{0}$$

- However, we can still evaluate $\nabla_{\mathbf{w}} \log P(\mathbb{X} | \mathbf{w})$ and use the iterative methods to solve \mathbf{w}
 - E.g., stochastic gradient descent

ML Estimation

- Log-likelihood:

$$\begin{aligned}\log P(\mathbb{X} | \mathbf{w}) &= \log \prod_{i=1}^N P(\mathbf{x}^{(i)}, y^{(i)} | \mathbf{w}) \\ &= \log \prod_i P(y^{(i)} | \mathbf{x}^{(i)}, \mathbf{w}) P(\mathbf{x}^{(i)} | \mathbf{w}) \\ &\propto \log \prod_i \sigma(\mathbf{w}^\top \mathbf{x}^{(i)})^{y^{(i)}} [1 - \sigma(\mathbf{w}^\top \mathbf{x}^{(i)})]^{(1-y^{(i)})} \\ &= \sum_i y^{(i)} \mathbf{w}^\top \mathbf{x}^{(i)} - \log(1 + e^{\mathbf{w}^\top \mathbf{x}^{(i)}}) \quad [\text{Homework}]\end{aligned}$$

- Unlike in linear regression, we cannot solve \mathbf{w} analytically in a closed form via

$$\nabla_{\mathbf{w}} \log P(\mathbb{X} | \mathbf{w}) = \sum_{t=1}^N [y^{(t)} - \sigma(\mathbf{w}^\top \mathbf{x}^{(t)})] \mathbf{x}^{(t)} = \mathbf{0}$$

- However, we can still evaluate $\nabla_{\mathbf{w}} \log P(\mathbb{X} | \mathbf{w})$ and use the iterative methods to solve \mathbf{w}
 - E.g., stochastic gradient descent
- It can be shown that $\log P(\mathbb{X} | \mathbf{w})$ is concave in terms of \mathbf{w} [1]
 - So, iterative algorithms converges

Outline

- 1 Probabilistic Models
- 2 Maximum Likelihood Estimation
 - Linear Regression
 - Logistic Regression
- 3 Maximum A Posteriori Estimation
- 4 Bayesian Estimation**

MAP Estimation

- So far, we solve \mathbf{w} by ML estimation:

$$\arg \max_{\mathbf{w}} P(\mathbb{X} | \mathbf{w})$$

MAP Estimation

- So far, we solve \mathbf{w} by ML estimation:

$$\arg \max_{\mathbf{w}} P(\mathbb{X} | \mathbf{w})$$

- In MAP estimation, we solve

$$\arg \max_{\mathbf{w}} P(\mathbf{w} | \mathbb{X}) = \arg \max_{\mathbf{w}} P(\mathbb{X} | \mathbf{w}) P(\mathbf{w})$$

- $P(\mathbf{w})$ models our *preference* or *prior knowledge* about \mathbf{w}

MAP Estimation for Linear Regression

- MAP estimation in linear regression:

$$\arg \max_{\mathbf{w}} \log[\mathbb{P}(\mathbb{X} | \mathbf{w})\mathbb{P}(\mathbf{w})]$$

MAP Estimation for Linear Regression

- MAP estimation in linear regression:

$$\arg \max_{\mathbf{w}} \log[\mathbb{P}(\mathbb{X} | \mathbf{w})\mathbb{P}(\mathbf{w})]$$

- If we assume that $\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \beta^{-1}\mathbf{I})$

$$\log[\mathbb{P}(\mathbb{X} | \mathbf{w})\mathbb{P}(\mathbf{w})] = \log\mathbb{P}(\mathbb{X} | \mathbf{w}) + \log\mathbb{P}(\mathbf{w})$$

MAP Estimation for Linear Regression

- MAP estimation in linear regression:

$$\arg \max_{\mathbf{w}} \log[\mathbf{P}(\mathbb{X} | \mathbf{w})\mathbf{P}(\mathbf{w})]$$

- If we assume that $\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \beta^{-1}\mathbf{I})$

$$\begin{aligned} \log[\mathbf{P}(\mathbb{X} | \mathbf{w})\mathbf{P}(\mathbf{w})] &= \log \mathbf{P}(\mathbb{X} | \mathbf{w}) + \log \mathbf{P}(\mathbf{w}) \propto -\sum_i (y^{(i)} - \mathbf{w}^\top \mathbf{x}^{(i)})^2 \\ &+ \log \sqrt{\frac{1}{(2\pi)^D \det(\beta^{-1}\mathbf{I})}} \exp \left[-\frac{1}{2} (\mathbf{w} - \mathbf{0})^\top (\beta^{-1}\mathbf{I})^{-1} (\mathbf{w} - \mathbf{0}) \right] \end{aligned}$$

MAP Estimation for Linear Regression

- MAP estimation in linear regression:

$$\arg \max_{\mathbf{w}} \log[\mathbb{P}(\mathbb{X} | \mathbf{w})\mathbb{P}(\mathbf{w})]$$

- If we assume that $\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \beta^{-1}\mathbf{I})$

$$\begin{aligned} \log[\mathbb{P}(\mathbb{X} | \mathbf{w})\mathbb{P}(\mathbf{w})] &= \log \mathbb{P}(\mathbb{X} | \mathbf{w}) + \log \mathbb{P}(\mathbf{w}) \propto -\sum_i (y^{(i)} - \mathbf{w}^\top \mathbf{x}^{(i)})^2 \\ &\quad + \log \sqrt{\frac{1}{(2\pi)^D \det(\beta^{-1}\mathbf{I})}} \exp \left[-\frac{1}{2} (\mathbf{w} - \mathbf{0})^\top (\beta^{-1}\mathbf{I})^{-1} (\mathbf{w} - \mathbf{0}) \right] \\ &\propto -\sum_i (y^{(i)} - \mathbf{w}^\top \mathbf{x}^{(i)})^2 - \beta \mathbf{w}^\top \mathbf{w} \end{aligned}$$

- $\mathbb{P}(\mathbf{w})$ corresponds to the *weight decay* term in Ridge regression

MAP Estimation for Linear Regression

- MAP estimation in linear regression:

$$\arg \max_{\mathbf{w}} \log[\mathbb{P}(\mathbb{X} | \mathbf{w})\mathbb{P}(\mathbf{w})]$$

- If we assume that $\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \beta^{-1}\mathbf{I})$

$$\begin{aligned} \log[\mathbb{P}(\mathbb{X} | \mathbf{w})\mathbb{P}(\mathbf{w})] &= \log \mathbb{P}(\mathbb{X} | \mathbf{w}) + \log \mathbb{P}(\mathbf{w}) \propto -\sum_i (y^{(i)} - \mathbf{w}^\top \mathbf{x}^{(i)})^2 \\ &\quad + \log \sqrt{\frac{1}{(2\pi)^D \det(\beta^{-1}\mathbf{I})}} \exp\left[-\frac{1}{2}(\mathbf{w} - \mathbf{0})^\top (\beta^{-1}\mathbf{I})^{-1}(\mathbf{w} - \mathbf{0})\right] \\ &\propto -\sum_i (y^{(i)} - \mathbf{w}^\top \mathbf{x}^{(i)})^2 - \beta \mathbf{w}^\top \mathbf{w} \end{aligned}$$

- $\mathbb{P}(\mathbf{w})$ corresponds to the **weight decay** term in Ridge regression
- MAP estimation provides a way to design complicated yet interpretable regularization terms
 - E.g., we have LASSO by letting $\mathbb{P}(\mathbf{w}) \sim \text{Laplace}(0, b)$ [Proof]
 - We can also let $\mathbb{P}(\mathbf{w})$ be a mixture of Gaussians

Remarks on ML and MAP Estimation

Theorem (Consistency)

The ML estimator Θ_{ML} is **consistent**, i.e., $\lim_{N \rightarrow \infty} \Theta_{ML} \xrightarrow{\text{Pr}} \Theta^*$ as long as the “true” $P(y|\mathbf{x}; \Theta^*)$ lies within our model \mathbb{F} .

Remarks on ML and MAP Estimation

Theorem (Consistency)

The ML estimator Θ_{ML} is **consistent**, i.e., $\lim_{N \rightarrow \infty} \Theta_{ML} \xrightarrow{\text{Pr}} \Theta^*$ as long as the “true” $P(y|x; \Theta^*)$ lies within our model \mathbb{F} .

Theorem (Cramér-Rao Lower Bound [2])

At a fixed (large) number N of examples, no consistent estimator of Θ^* has a lower expected MSE (mean square error) than the ML estimator Θ_{ML} .

- That is, Θ_{ML} has a low sample complexity (or is statistic efficient)

Remarks on ML and MAP Estimation

Theorem (Consistency)

The ML estimator Θ_{ML} is **consistent**, i.e., $\lim_{N \rightarrow \infty} \Theta_{ML} \xrightarrow{\text{Pr}} \Theta^*$ as long as the “true” $P(y|\mathbf{x}; \Theta^*)$ lies within our model \mathbb{F} .

Theorem (Cramér-Rao Lower Bound [2])

At a fixed (large) number N of examples, no consistent estimator of Θ^* has a lower expected MSE (mean square error) than the ML estimator Θ_{ML} .

- That is, Θ_{ML} has a low sample complexity (or is statistic efficient)
- ML estimation is popular due to its consistency and efficiency

Remarks on ML and MAP Estimation

Theorem (Consistency)

The ML estimator Θ_{ML} is **consistent**, i.e., $\lim_{N \rightarrow \infty} \Theta_{ML} \xrightarrow{\text{Pr}} \Theta^*$ as long as the “true” $P(y|\mathbf{x}; \Theta^*)$ lies within our model \mathbb{F} .

Theorem (Cramér-Rao Lower Bound [2])

At a fixed (large) number N of examples, no consistent estimator of Θ^* has a lower expected MSE (mean square error) than the ML estimator Θ_{ML} .

- That is, Θ_{ML} has a low sample complexity (or is statistic efficient)
- ML estimation is popular due to its consistency and efficiency
- When N is small that yields overfitting behavior, we can use MAP estimation to **introduce bias** and **reduce variance**

Remarks on ML and MAP Estimation

Theorem (Consistency)

The ML estimator Θ_{ML} is **consistent**, i.e., $\lim_{N \rightarrow \infty} \Theta_{ML} \xrightarrow{\text{Pr}} \Theta^*$ as long as the “true” $P(y|\mathbf{x}; \Theta^*)$ lies within our model \mathbb{F} .

Theorem (Cramér-Rao Lower Bound [2])

At a fixed (large) number N of examples, no consistent estimator of Θ^* has a lower expected MSE (mean square error) than the ML estimator Θ_{ML} .

- That is, Θ_{ML} has a low sample complexity (or is statistic efficient)
- ML estimation is popular due to its consistency and efficiency
- When N is small that yields overfitting behavior, we can use MAP estimation to **introduce bias** and **reduce variance**

Outline

- 1 Probabilistic Models
- 2 Maximum Likelihood Estimation
 - Linear Regression
 - Logistic Regression
- 3 Maximum A Posteriori Estimation
- 4 **Bayesian Estimation****

Bayesian Estimation

- In ML/MAP estimation, we solve Θ first, then uses it as a constant to make prediction:

$$\hat{y} = \arg \max_y P(y | \mathbf{x}; \Theta)$$

Bayesian Estimation

- In ML/MAP estimation, we solve Θ first, then uses it as a constant to make prediction:

$$\hat{y} = \arg \max_y P(y | \mathbf{x}; \Theta)$$

- *Bayesian estimation* treats Θ as a random variable:

$$\hat{y} = \arg \max_y P(y | \mathbf{x}, \mathbb{X}) = \arg \max_y \int P(y, \Theta | \mathbf{x}, \mathbb{X}) d\Theta$$

- Makes prediction by considering *all* Θ 's (weighted by their chances)

Bayesian Estimation

- In ML/MAP estimation, we solve Θ first, then uses it as a constant to make prediction:

$$\hat{y} = \arg \max_y P(y | \mathbf{x}; \Theta)$$

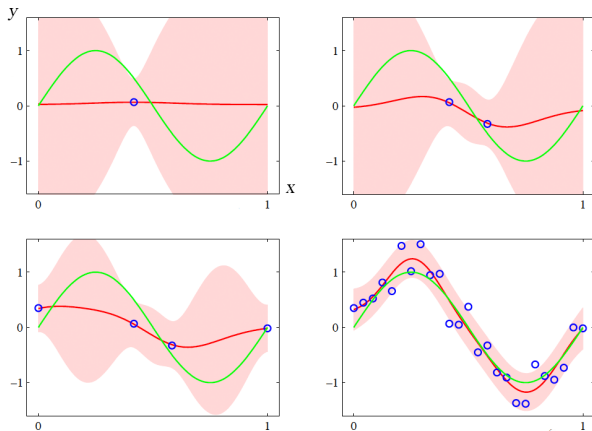
- *Bayesian estimation* treats Θ as a random variable:

$$\hat{y} = \arg \max_y P(y | \mathbf{x}, \mathbb{X}) = \arg \max_y \int P(y, \Theta | \mathbf{x}, \mathbb{X}) d\Theta$$

- Makes prediction by considering *all* Θ 's (weighted by their chances)
- Bayesian estimation usually generalizes much better when the size N of training set is small

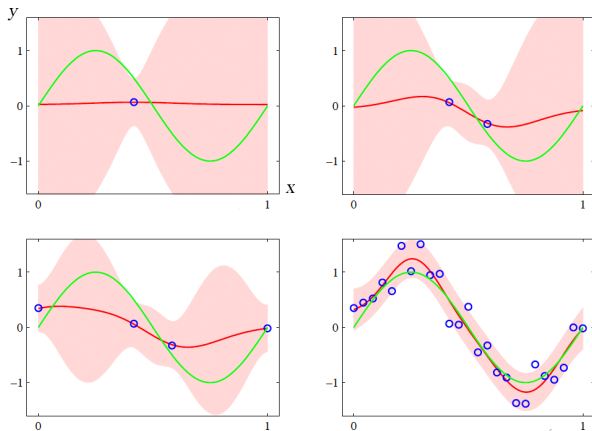
Bayesian vs. ML Estimation

- Example: polynomial regression



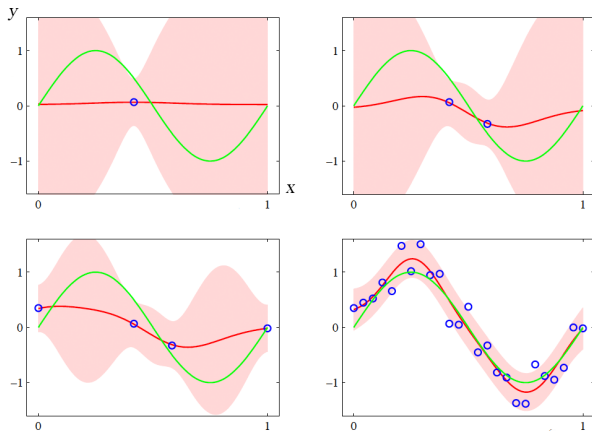
Bayesian vs. ML Estimation

- Example: polynomial regression
- Red line: predictions by Bayesian estimation regressor



Bayesian vs. ML Estimation

- Example: polynomial regression
- Red line: predictions by Bayesian estimation regressor
- Shaded area: predictions by ML/MAP estimation regressors

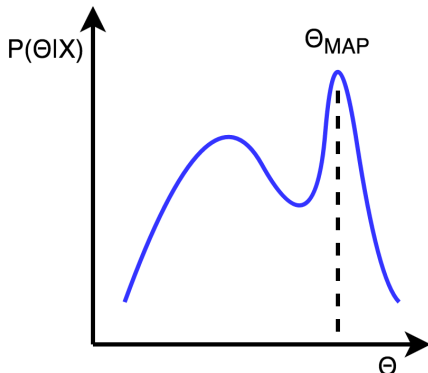


Bayesian vs. MAP Estimation

- MAP gains some benefit of Bayesian approach by incorporating prior as $\text{bias}(\Theta_{\text{MAP}})$
 - Reduces $\text{Var}_{\mathbb{X}}(\Theta_{\text{MAP}})$ when training set is small

Bayesian vs. MAP Estimation

- MAP gains some benefit of Bayesian approach by incorporating prior as bias(Θ_{MAP})
 - Reduces $\text{Var}_{\mathbb{X}}(\Theta_{\text{MAP}})$ when training set is small
- However, does **not** work if Θ_{MAP} is unrepresentative of the majority Θ in $\int P(\mathbf{y}, \Theta | \mathbf{x}, \mathbb{X}) d\Theta$
- E.g. when $P(\Theta | \mathbb{X})$ is a mixture of Gaussian



Remarks

- Bayesian estimation:

$$\hat{\mathbf{y}} = \arg \max_{\mathbf{y}} \mathbf{P}(\mathbf{y} | \mathbf{x}, \mathbb{X}) = \arg \max_{\mathbf{y}} \int \mathbf{P}(\mathbf{y}, \Theta | \mathbf{x}, \mathbb{X}) d\Theta$$

- Usually generalizes much better given a small training set

Remarks

- Bayesian estimation:

$$\hat{\mathbf{y}} = \arg \max_{\mathbf{y}} \mathbf{P}(\mathbf{y} | \mathbf{x}, \mathbb{X}) = \arg \max_{\mathbf{y}} \int \mathbf{P}(\mathbf{y}, \Theta | \mathbf{x}, \mathbb{X}) d\Theta$$

- Usually generalizes much better given a small training set
- Unfortunately, solution may not be tractable in many applications

Remarks

- Bayesian estimation:

$$\hat{\mathbf{y}} = \arg \max_{\mathbf{y}} \mathbf{P}(\mathbf{y} | \mathbf{x}, \mathbb{X}) = \arg \max_{\mathbf{y}} \int \mathbf{P}(\mathbf{y}, \Theta | \mathbf{x}, \mathbb{X}) d\Theta$$

- Usually generalizes much better given a small training set
- Unfortunately, solution may not be tractable in many applications
- Even tractable, incurs high computation cost
 - Not suitable for large-scale learning tasks

Reference I

- [1] Deepak Roy Chittajallu.
Why is the error function minimized in logistic regression convex?
<http://mathgotchas.blogspot.tw/2011/10/why-is-error-function-minimized-in.html>, 2011.
- [2] Harald Cramér.
Mathematical Methods of Statistics.
Princeton university press, 1946.