

IOWA STATE UNIVERSITY

Department of Computer Science

# COM S 573: Machine Learning

## Lecture 4: Probabilistic Interpretation of Linear Regression

# Linear Regression: Summary

- Representation

$$\hat{y}_i = \mathbf{w} \mathbf{x}_i^T$$

Predict a continuous scalar.

- Evaluation

$$RSS(\mathbf{w}) = (\mathbf{y} - \mathbf{X} \mathbf{w})^T (\mathbf{y} - \mathbf{X} \mathbf{w})$$

- Optimization

$$\frac{\partial RSS(\mathbf{w})}{\partial \mathbf{w}} = 0 \rightarrow \mathbf{w}^* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

# Linear Regression: Probabilistic interpretation

- Evaluation

$$RSS(\mathbf{w}) = (\mathbf{y} - \mathbf{X}\mathbf{w})^T (\mathbf{y} - \mathbf{X}\mathbf{w})$$

- Residual for each data sample

$$r_i = y_i - \mathbf{x}_i^T \mathbf{w}$$

What assumptions are we using?

# Linear Regression: Assumptions

- Data are linear

$$\hat{y}_i = \mathbf{w}^T \mathbf{x}_i$$

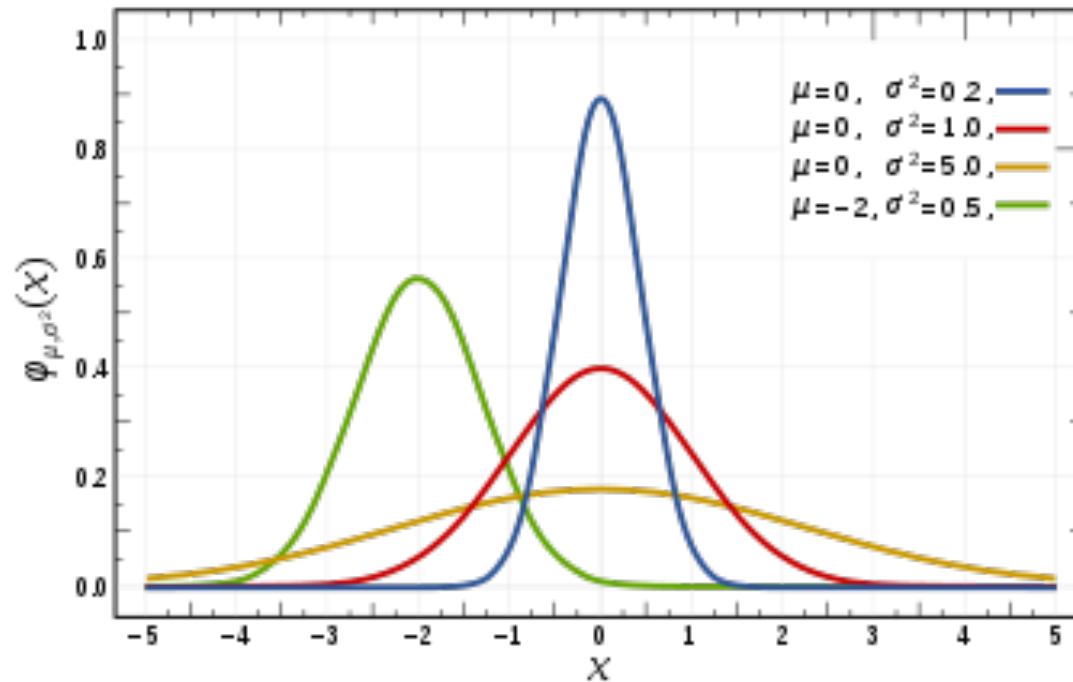
- Residuals are Independent and identically distributed

$$p(r_i, r_j) = p(r_i)p(r_j)$$

- Each residual follows normal distribution

$$r_i \sim \mathcal{N}(\mu, \sigma^2)$$

# Linear Regression: Normal Distribution



$$x \sim \mathcal{N}(\mu, \sigma^2) \rightarrow p(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

# Linear Regression: likelihood

- Given input samples

$$\mathbf{x} = \{x_1, x_2, \dots, x_n\}$$

- Likelihood vs. Probability

$$\mathcal{L}(\mathbf{w}|\mathbf{x}) \equiv p(\mathbf{x}|\mathbf{w})$$

- Likelihood of all input samples

$$\begin{aligned} p(\mathbf{x}|\mathbf{w}) &= p(x_1, x_2, \dots, x_n|\mathbf{w}) \\ &= p(x_1|\mathbf{w})p(x_2|\mathbf{w}) \cdots p(x_n|\mathbf{w}) \\ &= \prod_{i=1}^n p(x_i|\mathbf{w}) \end{aligned}$$

Why?



# Linear Regression: Optimization

- Maximum Likelihood estimator (MLE)

$$\boldsymbol{w}^* = \max_{\boldsymbol{w}} \mathcal{L}(\boldsymbol{w}|\boldsymbol{x}) = \max_{\boldsymbol{w}} \prod_{i=1}^n p(x_i|\boldsymbol{w})$$

- Log-likelihood

$$\log \mathcal{L}(\boldsymbol{w}|\boldsymbol{x}) = \log \prod_{i=1}^n p(x_i|\boldsymbol{w}) = \sum_{i=1}^n \log p(x_i|\boldsymbol{w})$$

# Linear Regression: Optimization

Why 0?



- Noisy observation model

$$y_i = \mathbf{w}^T \mathbf{x}_i + r_i, r_i \sim \mathcal{N}(0, \sigma^2)$$

- Probability

$$p(r) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(r-\mu)^2}{2\sigma^2}}$$

- Likelihood

$$\log \mathcal{L}(\mathbf{w}|\mathbf{x}) = \sum_{i=1}^n \log p(x_i|\mathbf{w})$$



# Linear Regression: Likelihood

- $\log \mathcal{L}(\mathbf{w}|\mathbf{x}) = \sum_{i=1}^n \log p(x_i|\mathbf{w})$

$$p(r) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(r-\mu)^2}{2\sigma^2}}$$

# Linear Regression: Likelihood

- $$\begin{aligned}\log \mathcal{L}(\mathbf{w}|\mathbf{x}) &= \sum_{i=1}^n \log p(x_i|\mathbf{w}) \\ &= \sum_{i=1}^n \log \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(\mathbf{y}_i - \mathbf{x}_i^T \mathbf{w})^2}{2\sigma^2}}\end{aligned}$$

$$p(r) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(r-\mu)^2}{2\sigma^2}}$$

# Linear Regression: Likelihood

- $$\begin{aligned}\log \mathcal{L}(\mathbf{w}|\mathbf{x}) &= \sum_{i=1}^n \log p(x_i|\mathbf{w}) \\ &= \sum_{i=1}^n \log \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(\mathbf{y}_i - \mathbf{x}_i^T \mathbf{w})^2}{2\sigma^2}} \\ &= \sum_{i=1}^n -\log \sqrt{2\pi}\sigma - \sum_{i=1}^n \frac{(\mathbf{y}_i - \mathbf{x}_i^T \mathbf{w})^2}{2\sigma^2}\end{aligned}$$

# Linear Regression: Likelihood

- $$\begin{aligned}\log \mathcal{L}(\mathbf{w}|\mathbf{x}) &= \sum_{i=1}^n \log p(x_i|\mathbf{w}) \\&= \sum_{i=1}^n \log \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(\mathbf{y}_i - \mathbf{x}_i^T \mathbf{w})^2}{2\sigma^2}} \\&= \sum_{i=1}^n -\log \sqrt{2\pi}\sigma - \sum_{i=1}^n \frac{(\mathbf{y}_i - \mathbf{x}_i^T \mathbf{w})^2}{2\sigma^2} \\&= -\frac{1}{2} \left( \frac{1}{\sigma^2} \sum_{i=1}^n (\mathbf{y}_i - \mathbf{x}_i^T \mathbf{w})^2 + n \log \sigma^2 + \text{const} \right)\end{aligned}$$

# Linear Regression: Likelihood

- $$\begin{aligned}\log \mathcal{L}(\mathbf{w}|\mathbf{x}) &= \sum_{i=1}^n \log p(x_i|\mathbf{w}) \\&= \sum_{i=1}^n \log \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(\mathbf{y}_i - \mathbf{x}_i^T \mathbf{w})^2}{2\sigma^2}} \\&= \sum_{i=1}^n -\log \sqrt{2\pi}\sigma - \sum_{i=1}^n \frac{(\mathbf{y}_i - \mathbf{x}_i^T \mathbf{w})^2}{2\sigma^2} \\&= -\frac{1}{2} \left( \frac{1}{\sigma^2} \sum_{i=1}^n (\mathbf{y}_i - \mathbf{x}_i^T \mathbf{w})^2 + n \log \sigma^2 + \text{const} \right) \\&= -\frac{1}{2} \left( \frac{1}{\sigma^2} \text{RSS}(\mathbf{w}) + n \log \sigma^2 \right) + \text{const}\end{aligned}$$

Maximize likelihood = minimize RSS

# Linear Regression: Likelihood

- Optimization

$$\log \mathcal{L}(\boldsymbol{w}|\boldsymbol{x}) = -\frac{1}{2} \left( \frac{1}{\sigma^2} RSS(\boldsymbol{w}) + n \log \sigma^2 \right) + \text{const}$$

- Find the optimal

$$\boldsymbol{w}^*, \sigma^*$$

HW1

# Linear Regression: Summary

- Representation

$$y_i = \mathbf{w}^T \mathbf{x}_i + r_i, r_i \sim \mathcal{N}(0, \sigma^2)$$

- Evaluation

$$\log \mathcal{L}(\mathbf{w}|\mathbf{x}) = \log \prod_{i=1}^n p(x_i|\mathbf{w}) = \sum_{i=1}^n \log p(x_i|\mathbf{w})$$

- Optimization

$$\log \mathcal{L}(\mathbf{w}|\mathbf{x}) = -\frac{1}{2} \left( \frac{1}{\sigma^2} RSS(\mathbf{w}) + n \log \sigma^2 \right) + \text{const}$$

# Linear Models: Next

- Linear regression: predict a scalar
  - House price
  - Weight of a planet
- Linear perceptron: classifier of discrete prediction
  - Predict an animal is a dog or not
  - Predict an image contains a square or not
- Logistic regression: classifier based on a probability
  - Predict how likely a team win
  - Predict how likely tomorrow is sunny