

A Review, then Sliding into Classification

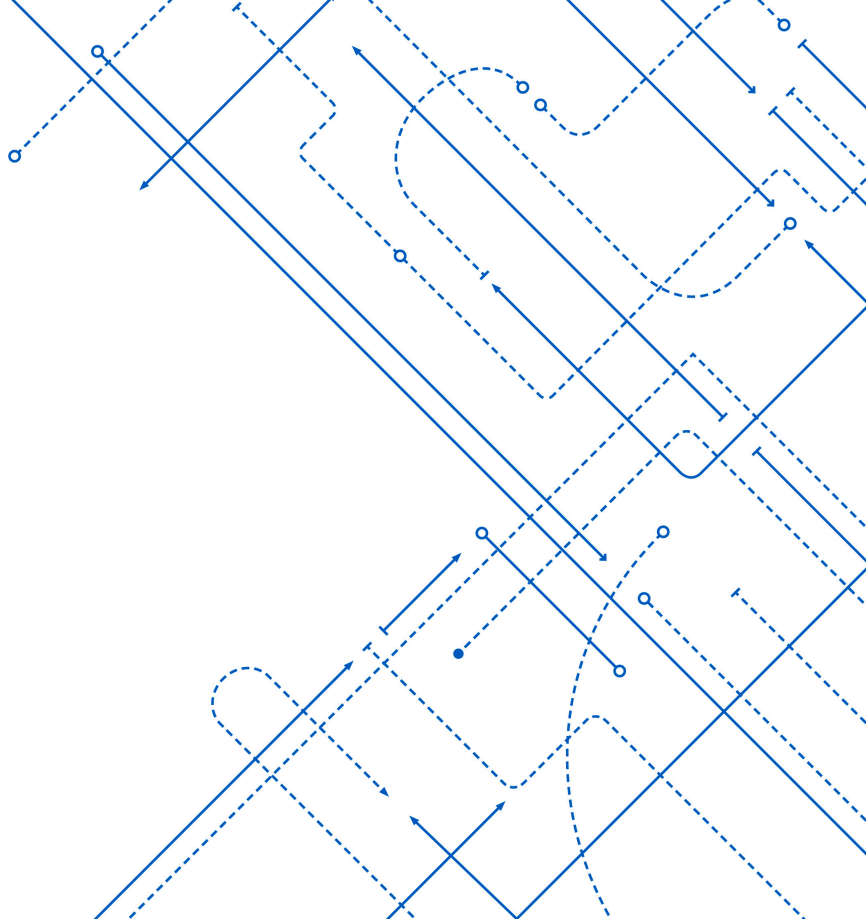
Kenneth (Kenny) Joseph



University at Buffalo

Department of Computer Science
and Engineering

School of Engineering and Applied Sciences

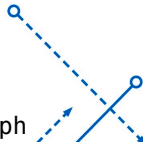


Announcements

- PA2 due Sunday night
- Quiz 4 is out, we will review Quiz 3 today
- PA 3 is out March 7th
 - Minimal coding, we'll be annotating data, calculating agreement statistics, and reflecting on the process
 - You have almost a month to do this (in other words, a break from programming assignments for a while)
- Midterm is March 17th
 - In class, mostly
 - One page handwritten notes, front and back
 - Official Accessibility requests **due by next Tuesday**
- Questions?

Terms/concepts you now know/have seen

- “Review”
 - Probability distribution
 - Expected Values
 - Stats (e.g. mean/variance)
 - Python
 - Pandas/Numpy/Jupyter
- ML High-level ideas
 - Model class
 - Loss function
 - Squared Error
 - Regularization
 - Optimization algorithm
 - (Stochastic) Gradient Descent
 - Closed form solutions
 - Making Predictions
- Models
 - (Regularized) Linear Regression
 - Polynomial regression
 - Decision Tree Regression
 - kNN regression
 - (Generalized) additive models
- Selection & Evaluation
 - 2/3-way holdout methods
 - K-fold cross validation
 - Bias/Variance tradeoff
 - Generalization error
 - The 3 sources of error
 - Over/underfitting



This week

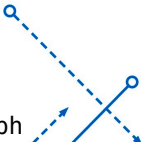
- PA 1, Quiz 3
- A brief review of where we're at
 - Supervised learning – what's the point?
 - Where do features come from?
 - What does `sklearn.linear_model.LinearRegression()` actually do?
- A “new” setup from a *probabilistic* perspective
 - Maximum Likelihood Estimation
 - Using the probabilistic approach to re-derive OLS regression
- Intro to classification
 - Logistic Regression
 - Bayes Optimal Classifier
 - Naïve Bayes
- Potentially: SVMs & Kernels

Back to the beginning

- In Supervised ML, we have...
 (x_i, y_i) where $x \in \mathbb{R}^d \leftarrow$ "Features" $y \leftarrow$ label
 $D = \{(x_1, y_1), \dots, (x_n, y_n)\} \in \mathbb{R}^d \times \mathcal{L}$

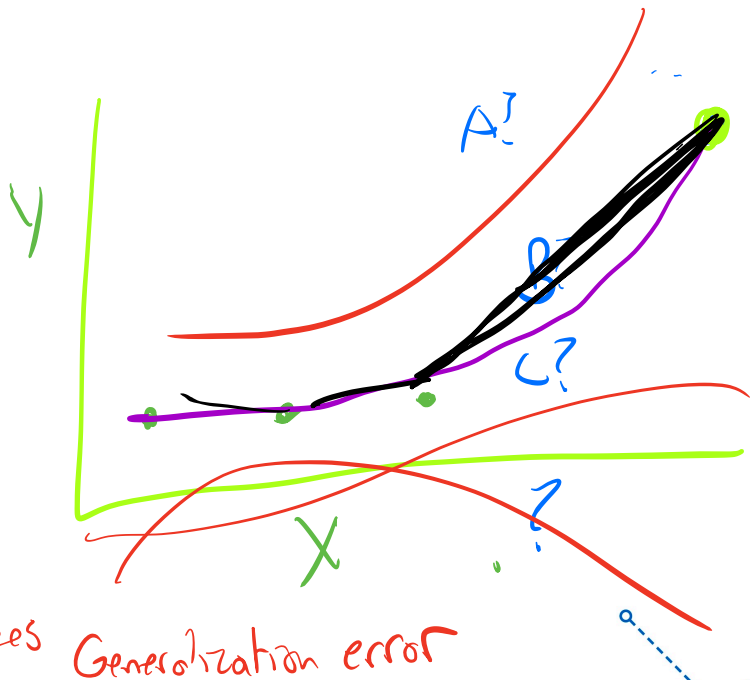
- We want to be able to get y when we only have $x...$

$$h(x_i) \rightarrow y_i$$



How do we find the “best model?”

- We could just memorize the training data ... right?
 - Training \neq Test, curse of dimensionality
- So... make assumptions (def'n model class)
- Then, find best model... 3 steps
 - Define best
 - Find best
 - Select/evaluate
- Linear Regression as an example...



How do we find the “best model?”

- We could just memorize the training data ... right?
 - Training \neq Test, curse of dimensionality
- So... make assumptions (def'n model class)
- Then, find best model... 3 steps
 - Define best
 - Find best
 - Select/evaluat
- Linear Regression as an example...

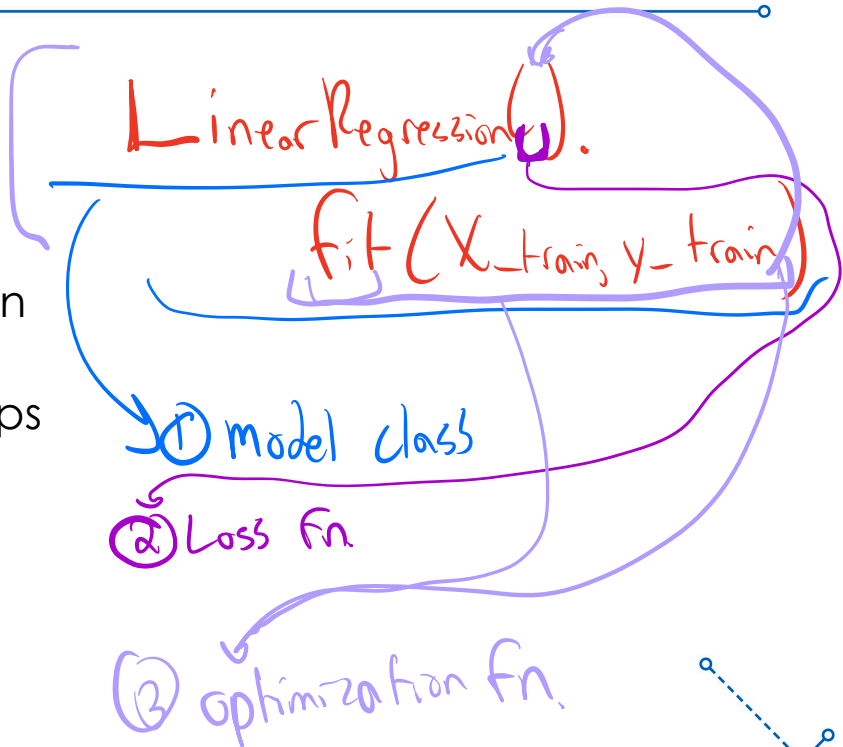
How do we find the "best model?"

- We could just memorize the training data ... right?
 - Training \neq Test, curse of dimensionality
- So... make assumptions (def'n model class)
- Then, find best model... 3 steps
 - Define best
 - Find best
 - Select/evaluate
- Linear Regression as an example...

LR
Model class? $f(x) = w^T x$
Loss fn? SSE $(w^T x - y)^2$
Optimize? GD / Closed Form
Selection/evaluation? RMSE

How do we find the “best model?”

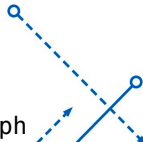
- We could just memorize the training data ... right?
 - Training \neq Test, curse of dimensionality
- So... make assumptions (def'n model class)
- Then, find best model... 3 steps
 - Define best
 - Find best
 - Select/evaluat
- Linear Regression as an example...



Re-introducing probability...

- Some holes in this “optimization” story...
 - What was all that business about “expectations”?
 - What about “training data as a random sample”? Of what?
 - Why SSE?
- Remembering the **probabilistic** part...

$$(x_i, y_i) \sim P(X, Y)$$



Implications of probabilistic framing

- Goal changes slightly – find $h(x)$ approx. y
- Re-specifying “the best we can do”...
- Re-explaining “training data as a random sample”...
- But what about the SSE optimization part?

$$\underbrace{h(x)} \approx y$$
$$E_{(x,y) \sim p} [L(x,y | h(\cdot))]$$

$$(x_i, y_i) \sim p$$

Maximum Likelihood estimation

- P(coin) is heads, when $D = \{H, T, T, H, H, H, T, T, T, T\}$?
- More formal derivation?
- Use **MLE**
 - Specify parameterized distribution
 - Find parameters that make observed data most likely
- For coin toss...

$$P(H) = \frac{n_H}{n_H + n_T} = 0.4$$
$$P(D|\theta) = \binom{n_H + n_T}{n_H} \theta^{n_H} (1-\theta)^{n_T}$$
$$\hat{\theta} = \operatorname{argmax}_{\theta} P(D; \theta)$$
$$\hat{\theta} = \operatorname{argmax}_{\theta}$$

Example from:

<https://www.cs.cornell.edu/courses/cs4780/2018fa/lectures/lecturenote06.html>

Linking back to the optimization view

- If we have a good model of the distrn. from which (\mathbf{x}, y) is drawn, we can use it to put forward a good guess as to $E[Y | X=x]$
- MLE gets us the best estimate of this probability distribution, **given a particular parameterized form...**

▪ Actually, two kinds of models

- Generative
- Discriminative

$(\mathbf{x}, y) \sim P$

$E_{P(\mathbf{x}, y)} [Y | X=x]$

$P(y|x)$

$P(\mathbf{x}, y)$

A "new strategy" for ML

1. Define $p(y | \mathbf{x})$ in terms of some parameterized model
2. Write down (log) likelihood function \leftarrow "loss fn."
3. Find the parameters that maximize the probability of the observed data

} "model class"

$$p(y_i | x_i) \sim N(w^T x_i, \sigma^2)$$

$$= \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x_i^T w - y_i)^2}{2\sigma^2}}$$

likelihood fn:

argmax_w

$$\prod p(y_i | x_i; w)$$

$$\begin{aligned} &\log \prod p(y_i | x_i; w) \\ &\in \log P(y, x; w) \end{aligned}$$



Trying our "new strategy" for linear regression

$$\operatorname{argmax}_w \sum \log p(y_i | x_i, w)$$

$$\sum \log \left(\frac{1}{\sqrt{2\pi}} \right) + \log(e^{-\dots})$$

See Sections 1 and 2 of <https://www.cs.cornell.edu/courses/cs4780/2018fa/lectures/lecturenote08.html> for the full derivation!

$$\sum \log(e^{-\frac{1}{2}(x_i^T w - y_i)^2})$$

$$\operatorname{argmax}_w - \sum_i (x_i^T w - y_i)^2$$

$$\operatorname{argmin}_w \sum (x_i^T w - y_i)^2$$

