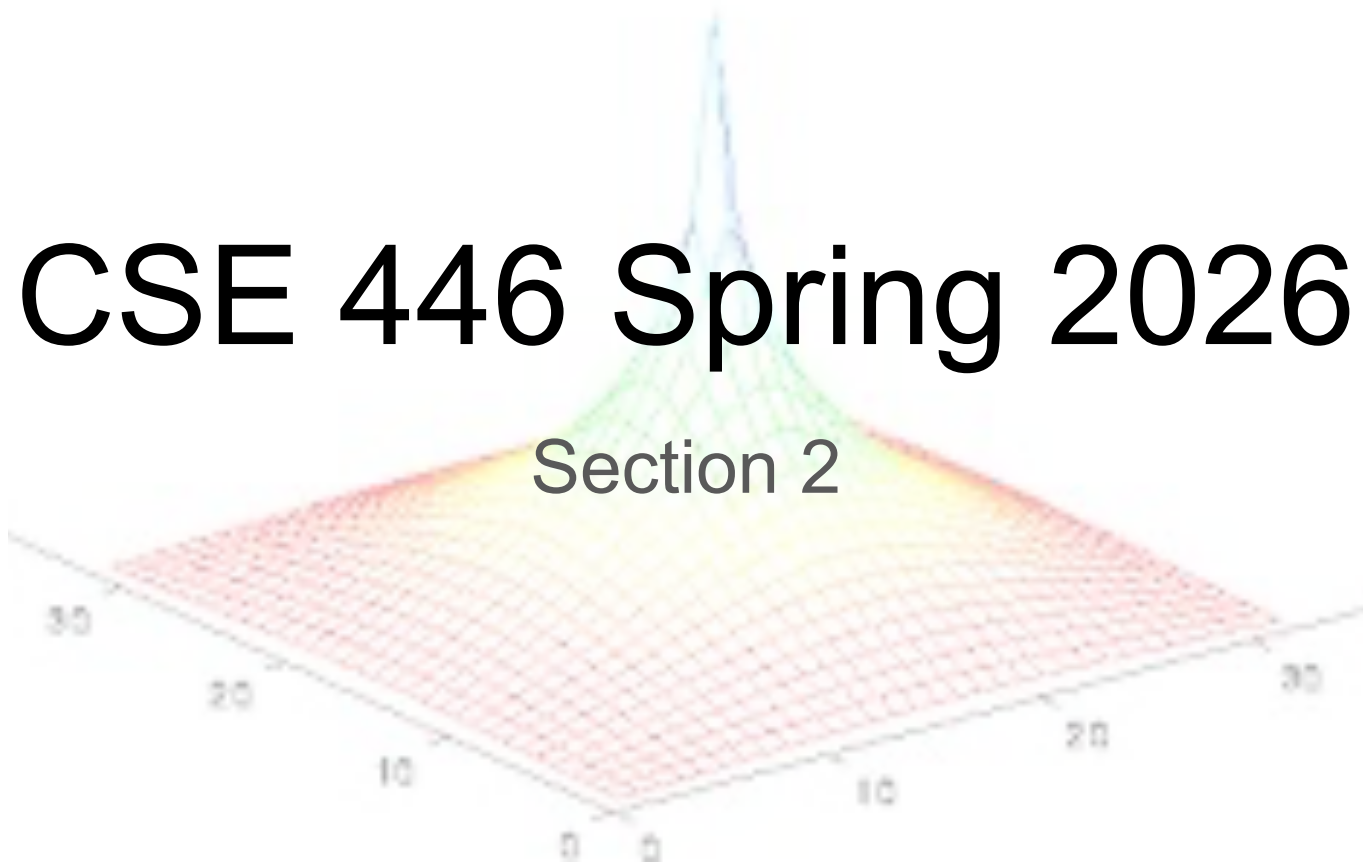


CSE 446 Spring 2026

Section 2



Reminders

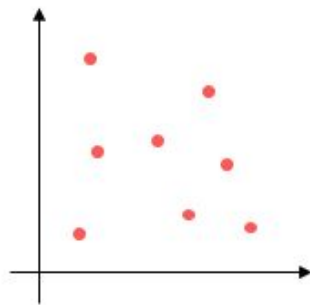
- **HW1 due Wed, Apr 22**

Some tips:

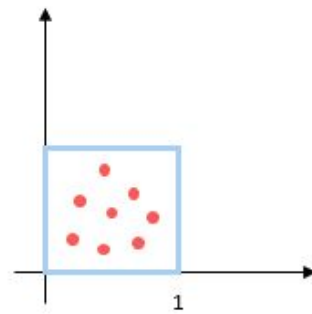
- Use office hours to your advantage
 - Student TA OH for homework questions
 - Professor OH for conceptual questions
- Skim the homeworks the day they are assigned and try one problem
 - Motivates you to get things done on time, starting an untouched assignment can be daunting
- Keep a tab open with the lecture slides while you do the homework for reference

- a. Data Normalization/Standardization (1.1/1.2)
- b. Train/Val/Test split (2.1)
- c. MLE (3.b)
- d. Review Linear Regression

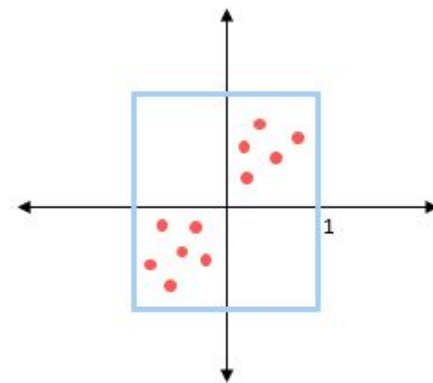
Problem 1



Actual Data

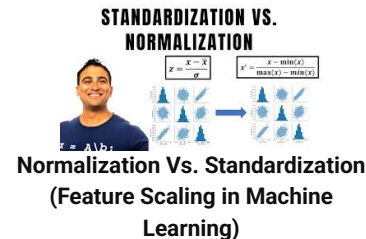


After normalizing



After standardization

Example of Normalization/Standardization



Stock Prediction Dataset

RAW ORIGINAL DATASET

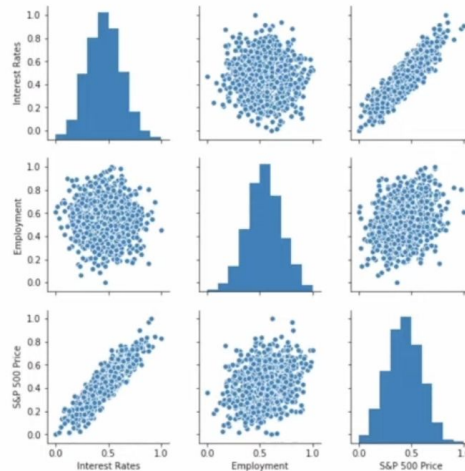
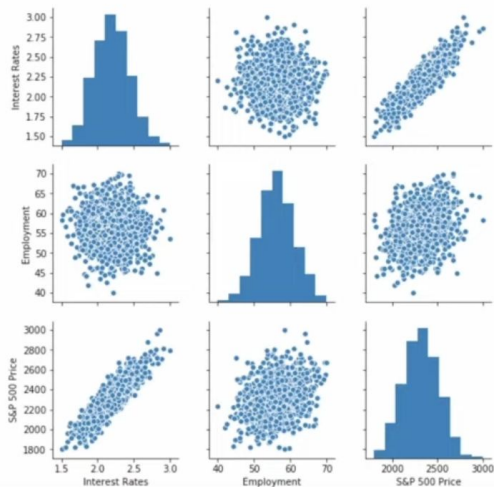
	Interest Rates	Employment	S&P 500 Price
0	1.943859	55.413571	2206.680582
1	2.258229	59.546305	2486.474488
2	2.215863	57.414687	2405.868337
3	1.977960	49.908353	2140.434475
4	2.437723	52.035492	2411.275663
5	2.143637	56.060598	2187.344909
6	2.148647	51.513208	2263.049249
7	2.176184	53.475909	2281.496374
8	2.125352	63.668422	2355.163011
9	2.225682	56.993396	2326.330337

QUICK STATS!

	Interest Rates	Employment	S&P 500 Price
count	1000.00	1000.00	1000.00
mean	2.20	56.25	2320.00
std	0.24	4.86	193.85
min	1.50	40.00	1800.00
25%	2.04	53.04	2190.45
50%	2.20	56.16	2312.44
75%	2.36	59.42	2455.76
max	3.00	70.00	3000.00

Normalization: Range 0 - 1

$$\tilde{x}_i^{(j)} = \frac{x_i^{(j)} - x_i^{\min}}{x_i^{\max} - x_i^{\min}}$$



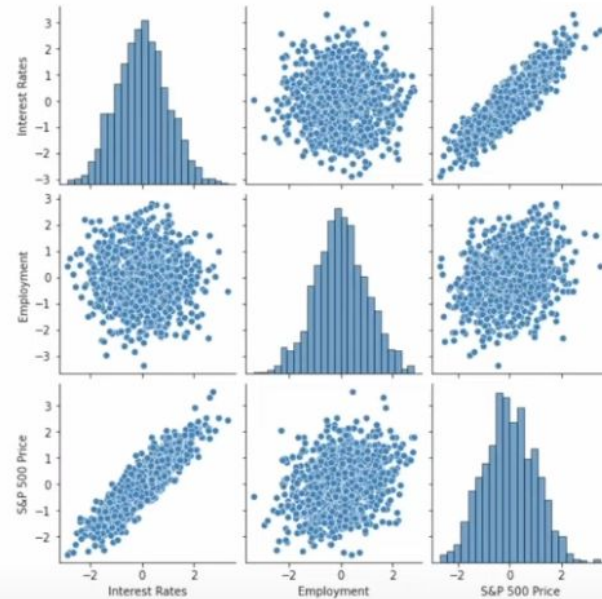
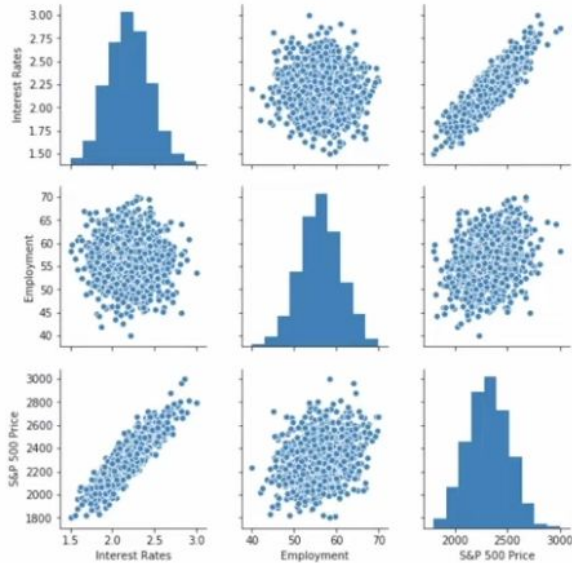
	Interest Rates	Employment	S&P 500 Price
count	1000.00	1000.00	1000.00
mean	2.20	56.25	2320.00
std	0.24	4.86	193.85
min	1.50	40.00	1800.00
25%	2.04	53.03	2190.45
50%	2.20	56.16	2312.44
75%	2.36	59.42	2455.76
max	3.00	70.00	3000.00



	Interest Rates	Employment	S&P 500 Price
count	1000.00	1000.00	1000.00
mean	0.46	0.54	0.43
std	0.16	0.16	0.16
min	0.00	0.00	0.00
25%	0.36	0.43	0.33
50%	0.47	0.54	0.43
75%	0.57	0.65	0.55
max	1.00	1.00	1.00

Standardization: 0 Mean, 1 Standard Deviation

$$\tilde{x}_i^{(j)} = \frac{x_i^{(j)} - \mu_i}{\sigma_i}$$



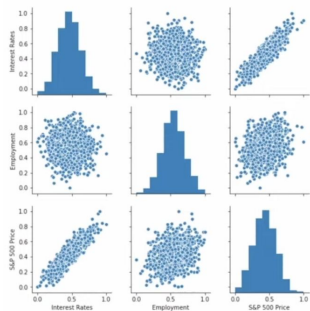
	Interest Rates	Employment	S&P 500 Price
count	1000.00	1000.00	1000.00
mean	2.20	56.25	2320.00
std	0.24	4.86	193.85
min	1.50	40.00	1800.00



	Interest Rates	Employment	S&P 500 Price
count	1000.00	1000.00	1000.00
mean	0.00	0.00	-0.00
std	1.00	1.00	1.00
min	-2.88	-3.34	-2.68

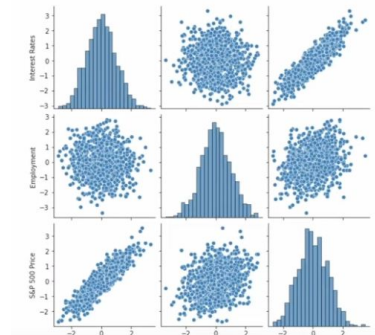
STANDARDIZATION Vs. NORMALIZATION?

- Generally speaking, there is no right or wrong answer!
- In case of neural networks, normalization is preferred since we don't assume any data distribution.
- Standardization is preferred when data follows gaussian distribution
- Standardization is preferred over normalization when there are a lot of outliers.



	Interest Rates	Employment	S&P 500 Price
count	1000.00	1000.00	1000.00
mean	0.46	0.54	0.43
std	0.16	0.16	0.16
min	0.00	0.00	0.00
25%	0.36	0.43	0.33
50%	0.47	0.54	0.43
75%	0.57	0.65	0.55
max	1.00	1.00	1.00


Task/Model Type	Prefer...	Why?
Neural Networks	Normalization	Helps convergence, bounded inputs
KNN, K-Means	Normalization	Distance-based, scale matters
Linear Models (LR, SVM)	Standardization	Assumes normally distributed features
PCA / LDA	Standardization	Sensitive to variance



	Interest Rates	Employment	S&P 500 Price
count	1000.00	1000.00	1000.00
mean	0.00	0.00	-0.00
std	1.00	1.00	1.00
min	-2.88	-3.34	-2.68

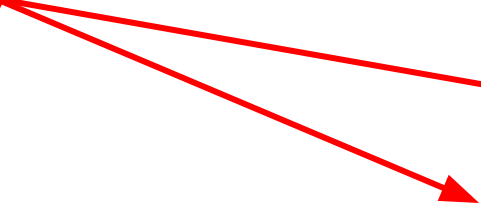
What do you never ever ever ever ever ever ever ever do?

Tune your model on your test set!



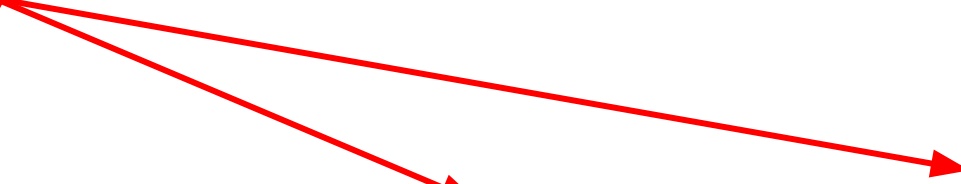
Training
parameters

Easy to
remember...



Choosing
hyperparameters

Easy-ish to
remember...



Preprocessing, small
changes, etc...

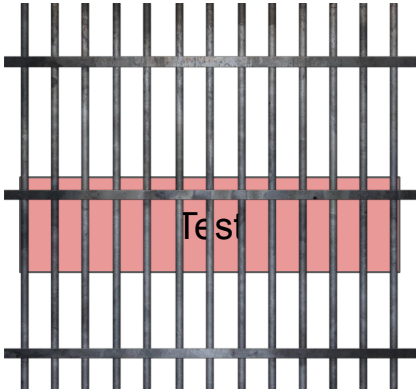
**Hard to
remember**
Very sinister!

Easiest way to combat this?

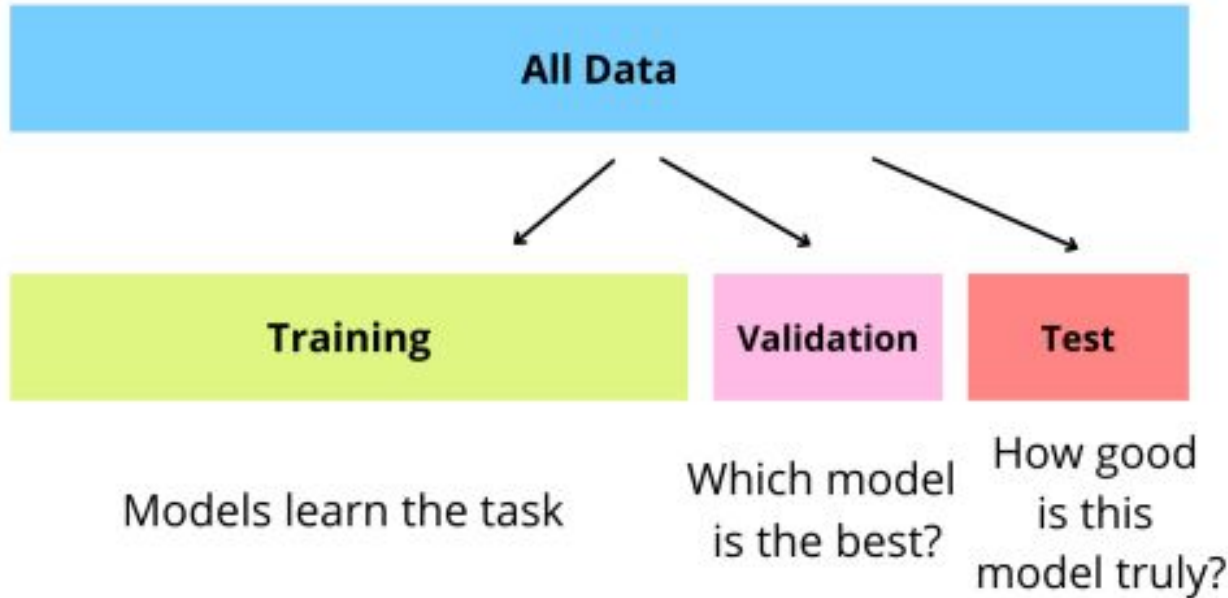


Only touch this until
everything is done

Then you
can free
him



Train/Val/Test Splits



Stock Prediction Dataset (E.g. train/val split)

	$X_1^{(1)}$	$X_1^{(2)}$	Y	
	Interest Rates	Employment	S&P 500 Price	
Training	0	1.943859	55.413571	2206.680582
	1	2.258229	59.546305	2486.474488
	2	2.215863	57.414687	2405.868337
	3	1.977960	49.908353	2140.434475
	4	2.437723	52.035492	2411.275663
	5	2.143637	56.060598	2187.344909
Validation	6	2.148647	51.513208	
	7	2.176184	53.475909	
	8	2.125352	63.668422	
	9	2.225682	56.993396	

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

Problem 1.1

Data Standardization

Data standardization is the task of transforming each feature in our dataset to have mean 0 and variance 1. The typical way to do this is using the Z-Score, which is defined as below:

$$\tilde{x}_i^{(j)} = \frac{x_i^{(j)} - \mu_i}{\sigma_i}$$

Where μ_i is the mean of each feature and σ_i is the standard deviation of each feature, which are empirically calculated from the data.

Question: what should you do when $\sigma_i = 0$ for some i ?

Constant feature.

Likely remove the feature before standardization and add it back to the dataset to avoid divide by zero error.

When training and evaluating your model, you should calculate the parameters for your normalization or standardization function on the training set **ONLY!**

In other words, if we were using Z-Score, we'd calculate our μ_i, σ_i on the training set, and use those same values when standardizing our validation/test data. The same applies to normalization methods, such as min-max scaling, where we want to determine x_i^{min}, x_i^{max} from training data. This will likely mean that we may get some values outside $[0, 1]$ after rescaling new data, but, given that our training dataset is large enough, this shouldn't be much of an issue (so unseen data likely won't be wildly outside of $[0, 1]$).

Question 1: Should we always choose x_i^{min} and x_i^{max} based on train data? Can we sometimes do better? Think about cases when we have some underlying information about data.

Consider RGB images. These are typically encoded as arrays of shape (3, height, width), with each value being an integer in range $[0, 255]$. In this case we should just use $x_i^{max} = 255$ and $x_i^{min} = 0$ to normalize the data.

Question 2: When can values outside of $[0, 1]$ range in test set cause issues?

Consider a model $f(x) = \log(x)^T w$. In this case if test datapoint have a value below 0, this code will fail, as log has domain $[0, \infty)$.

Problem 2.1

Program 1

```
1 # Given dataset of 1000-by-50 feature
2 # matrix X, and 1000-by-1 labels vector
3
4 mu = np.mean(X, axis=0)
5 X = X - mu
6
7 idx = np.random.permutation(1000)
8 TRAIN = idx[0:900]
9 TEST = idx[900::]
10
11 ytrain = y[TRAIN]
12 Xtrain = X[TRAIN, :]
13
14 # solve for argmin_w ||Xtrain*w - ytrain||_2
15 w = np.linalg.solve(np.dot(Xtrain.T, Xtrain), np.dot(Xtrain.T, ytrain))
16
17 b = np.mean(ytrain)
18
19 ytest = y[TEST]
20 Xtest = X[TEST, :]
21
22 train_error = np.dot(np.dot(Xtrain, w)+b - ytrain,
23                     np.dot(Xtrain, w)+b - ytrain ) / len(TRAIN)
24 test_error = np.dot(np.dot(Xtest, w)+b - ytest,
25                    np.dot(Xtest, w)+b - ytest ) / len(TEST)
26
27 print('Train error = ', train_error)
28 print('Test error = ', test_error)
```

μ is calculated from the **entire** data (train + test), intertwining them!

This is bad!

Correct procedure:

- **Split into train and test**
- **Compute the mean of the train data (μ_{train})**
- **De-mean both train and test data using μ_{train}**

Program 2

```
1 # We are given: 1) dataset X with n=1000 samples and 50 features and 2) a vector y of length 1000 with labels.
2 # Consider the following code to train a model, using cross validation to perform hyperparameter tuning.
3
4 def fit(Xin, Yin, _lambda):
5     w = np.linalg.solve(np.dot(Xin.T, Xin) + _lambda * np.eye(Xin.shape[1]), np.dot(Xin.T, Yin))
6     b = np.mean(Yin) - np.dot(w, mu)
7     return w, b
8
9 def predict(w, b, Xin):
10    return np.dot(Xin, w) + b
11
12 idx = np.random.permutation(1000)
13 TRAIN = idx[0:800]
14 VAL = idx[800:900]
15 TEST = idx[900::]
16
17 ytrain = y[TRAIN]
18 Xtrain = X[TRAIN, :]
19 yval = y[VAL]
20 Xval = X[VAL, :]
21
22 # demean data
23 mu = np.mean(Xtrain, axis=0)
24 Xtrain = Xtrain - mu
25 Xval = Xval - mu
26
27 # use validation set to pick the best hyper-parameter to use
28 lambdas = [10 ** -5, 10 ** -4, 10 ** -3, 10 ** -2]
29 err = np.zeros(len(lambdas))
30
31 for idx, _lambda in enumerate(lambdas):
32     w, b = fit(Xtrain, ytrain, _lambda)
33     yval_hat = predict(w, b, Xval)
34     err[idx] = np.mean((yval_hat - yval)**2)
35
36 lambda_best = lambdas[np.argmin(err)]
37
38 Xtot = np.concatenate((Xtrain, Xval), axis=0)
39 ytot = np.concatenate((ytrain, yval), axis=0)
40
41 w, b = fit(Xtot, ytot, lambda_best)
42
43 ytest = y[TEST]
44 Xtest = X[TEST, :]
45
46 # demean data
47 Xtest = Xtest - mu
48
49 ytot_hat = predict(w, b, Xtot)
50 train_error = np.mean((ytot_hat - ytot) **2)
51 ytest_hat = predict(w, b, Xtest)
52 test_error = np.mean((ytest_hat - ytest) **2)
53
54 print('Train error = ', train_error)
55 print('Test error = ', test_error)
```

The final model is trained on BOTH the training and validation sets.

This is... eh...

Your hyperparameters selected on just the train data may not hold for train + val:

- More data is good but you should ensure that the hyperparameters you tuned do not depend on the number of elements.
- Tradeoff between more data and better test error estimate

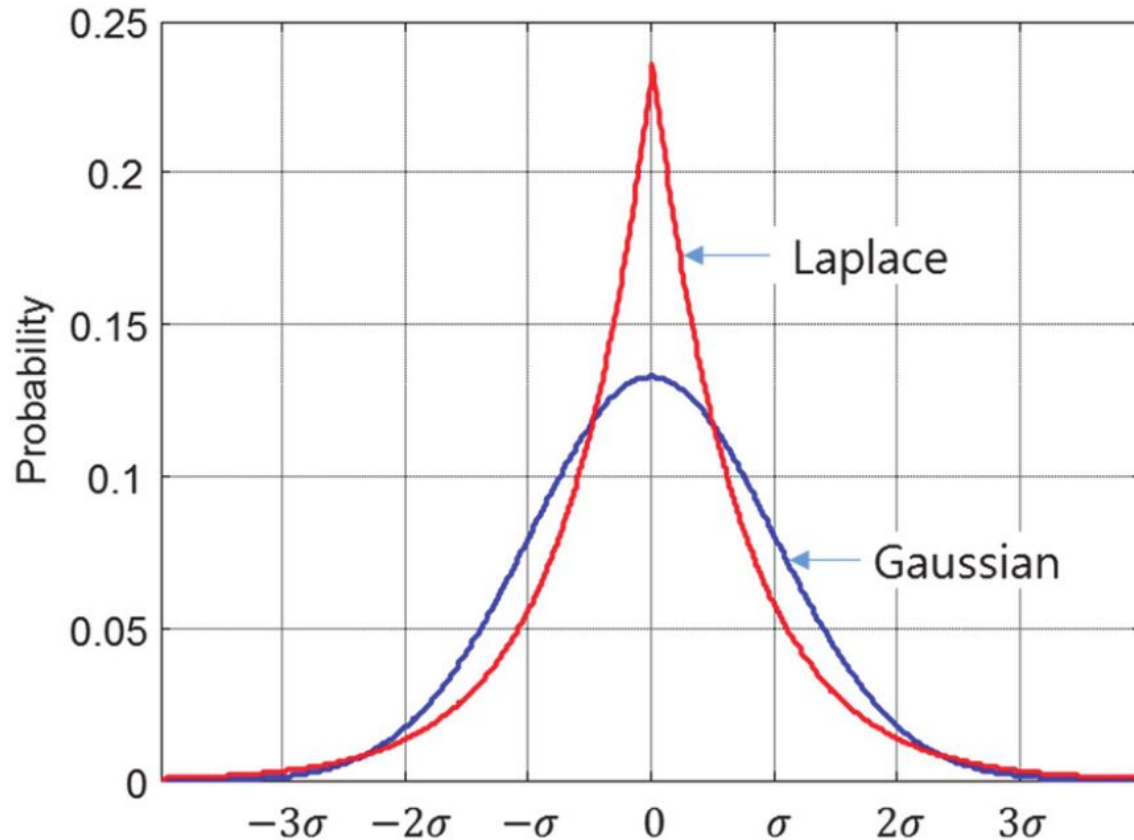
Problem 3b

Estimating from noisy measurements



- E.g. If your hand is not perfectly still, you introduce some noise
- The noise follows some distribution if studied
 - [e.g. more likely that your hand drops lower than goes higher]

Two Common Noise Distributions In The Class



Problem 3b - Context

Linear model with noise: $y_i = x_i^T w + v_i$

Rewrite as: $v_i = y_i - x_i^T w$

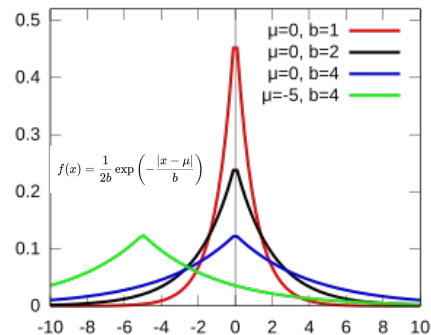
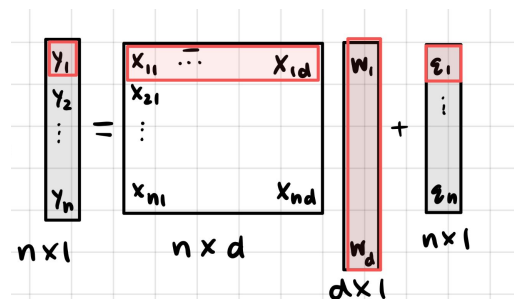
Noise has density $p(z=v)$: $p(z) = (1/2a) \exp(-|z|/a)$

Then:

$$\log p_w((x_1, y_1), \dots, (x_m, y_m)) = \log \prod_{i=1}^m p(y_i - x_i^T w)$$

ie. the probability of observing the data

We will use MLE to maximize the likelihood of seeing the data by finding \hat{w}_{MLE}



$\mu = 0$ and $b = a$

(which color curve best matches the noisy data distribution?)

Maximum Likelihood Estimation (MLE)

$$P(w | X, Y) = \frac{P(X, Y | w) P(w)}{P(X, Y)}$$

$$\text{posterior} = \frac{\text{likelihood} \cdot \text{prior}}{\text{evidence}}$$

$$\hat{w}_{MLE} = \arg \max_w p(X, Y | w)$$

Suppose we have data $\mathcal{D} = \{x^{(i)}\}_{i=1}^N$

$$\theta^{MLE} = \arg \max_{\theta} \prod_{i=1}^N p(\mathbf{x}^{(i)} | \theta)$$

Maximum Likelihood
Estimate (MLE)

$$\hat{w}_{MLE} = \arg \max_w \log p_w((x_1, y_1), \dots, (x_m, y_m))$$

$$\theta^{MLE} = \operatorname{argmax}_{\theta} \prod_{i=1}^N p(\mathbf{x}^{(i)} | \theta)$$

$$= \arg \max_w \log \prod_{i=1}^m p(y_i - x_i^T w)$$

	$X_1^{(1)}$ Interest Rates	$X_1^{(2)}$ Employment	Y S&P 500 Price
0	1.943859	55.413571	2206.680582
1	2.258229	59.546305	2486.474488
2	2.215863	57.414687	2405.868337
3	1.977960	49.908353	2140.434475
4	2.437723	52.035492	2411.275663
5	2.143637	56.060598	2187.344909

Let's look at the steps *before* considering the distribution:

3b) Laplacian - Work in pairs for 3 minutes for the first step

- (b) When the noise measurements follow a Laplacian distribution ($p(z) = (1/2a) \exp(-|z|/a)$), what is the maximum likelihood estimate of w ? Express your answer as the solution to an optimization problem such as in the previous part.

$$\hat{w}_{MLE} = \arg \max_w \log p_w((x_1, y_1), \dots, (x_m, y_m))$$

$$= \arg \max_w \log \prod_{i=1}^m p(y_i - x_i^\top w)$$

	x_1	x_2	y
0	1.943850	95.413071	2206.890582
1	2.258229	98.548305	2408.471448
2	2.215863	87.414987	2405.890307
3	1.977680	88.990353	2148.024475
4	2.437723	92.036482	2411.275983
5	2.143637	96.000598	2187.344009

$$= \arg \max_w \sum_{i=1}^m \log \underbrace{p(y_i - x_i^\top w)}_{\text{What do we plug in here?}} \quad [\log(ab) = \log a + \log b]$$

$$= \arg \max_w \sum_{i=1}^m \log \left(\frac{1}{2a} \cdot e^{-\frac{|y_i - x_i^T w|}{a}} \right)$$

$$= \arg \max_w \sum_{i=1}^m \log \left(\frac{1}{2a} \right) - \frac{|y_i - x_i^T w|}{a}$$

$$= \arg \max_w \sum_{i=1}^m -\frac{|y_i - x_i^T w|}{a}$$

constant offset doesn't affect opti

$$= \frac{1}{a} \cdot \arg \max_w \sum_{i=1}^m -|y_i - x_i^T w|$$

$$= \arg \max_w \sum_{i=1}^m -|y_i - x_i^T w|$$

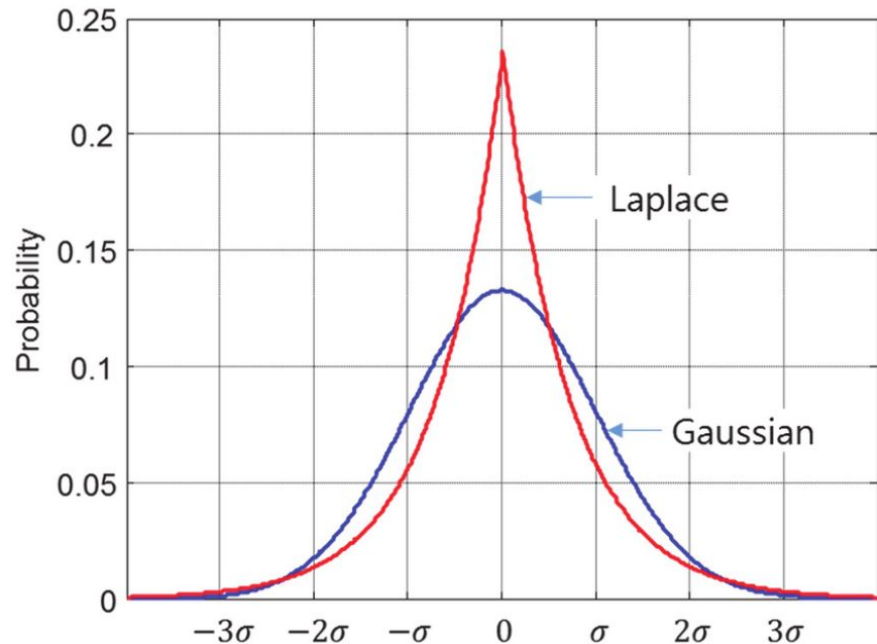
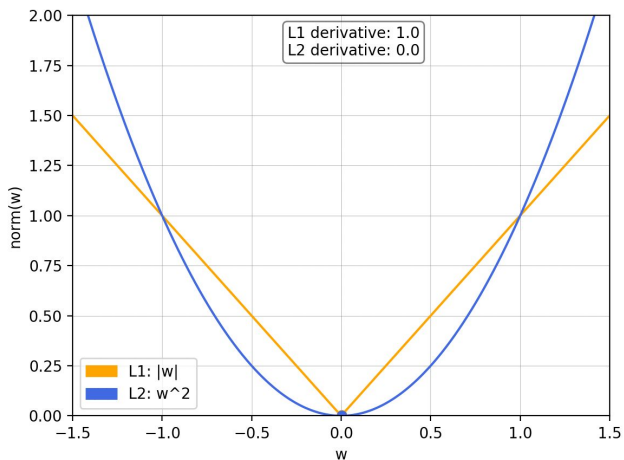
constant offset doesn't affect opti

$$= \arg \min_w \sum_{i=1}^m |y_i - x_i^T w| = \arg \min_w \|Xw - Y\|_1$$

Therefore the maximum likelihood estimate of w is $\hat{w} = \arg \min_w \|Xw - Y\|_1$.

Cool Connection

Visual similarities exist between the L1 norm vs. the Laplacian, and the L2 norm vs. the Gaussian



Using Laplacian (sharp) as noise results in the L1 norm (sharp) in the optimization equation

Using Gaussian (smooth) as noise results in the L2 norm (smooth) in the optimization equation

Linear Regression Review

We will continue the optimization problem for the standard linear regression (gaussian noise assumption)

Linear Regression

$$\hat{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmin}} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

Find \mathbf{w} that minimizes the SSE. This is $\hat{\mathbf{w}}$.

$$= \underset{\mathbf{w}}{\operatorname{argmin}} \sum_{i=1}^n (y_i - (\mathbf{x}_i^T \mathbf{w} + b))^2$$

Plug in $\hat{y}_i = \mathbf{x}_i^T \mathbf{w} + b$.

$$= \underset{\mathbf{w}}{\operatorname{argmin}} \sum_{i=1}^n (y_i - (\mathbf{x}_i^T \mathbf{w}))^2$$

Disregard b , since an intercept can be represented by appending a 1 to \mathbf{x} .

$$0 = \frac{\partial}{\partial \mathbf{w}} \sum_{i=1}^n (y_i - (\mathbf{x}_i^T \hat{\mathbf{w}}))^2$$

Find $\hat{\mathbf{w}}$ by taking the partial derivative of the argmin term and setting it equal to 0.

$$= \sum_{i=1}^n \frac{\partial}{\partial \mathbf{w}} (y_i - (\mathbf{x}_i^T \hat{\mathbf{w}}))^2$$

Property of derivatives.

$$= \sum_{i=1}^n 2(y_i - \mathbf{x}_i^T \hat{\mathbf{w}})(-\mathbf{x}_i)$$

Derive.

Linear Regression

$$= \sum_{i=1}^n (y_i - \mathbf{x}_i^T \hat{\mathbf{w}}) \mathbf{x}_i$$

Divide by -2

$$= \sum_{i=1}^n \mathbf{x}_i (y_i - \mathbf{x}_i^T \hat{\mathbf{w}})$$

The term $y_i - \mathbf{x}_i^T \hat{\mathbf{w}}$ is a scalar. $c\mathbf{x} = \mathbf{x}c$

$$= \sum_{i=1}^n (\mathbf{x}_i y_i - \mathbf{x}_i \mathbf{x}_i^T \hat{\mathbf{w}})$$

Distribute

$$= \sum_{i=1}^n \mathbf{x}_i y_i - \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T \hat{\mathbf{w}}$$

Distribute

$$\left(\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T \right) \hat{\mathbf{w}} = \sum_{i=1}^n \mathbf{x}_i y_i$$

Move second term to LHS.

$$\left(\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T\right) \hat{\mathbf{w}} = \sum_{i=1}^n \mathbf{x}_i y_i$$

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

$$\begin{aligned} (X^T X)^{-1} (X^T X) \hat{\mathbf{w}} &= (X^T X)^{-1} X^T \mathbf{y} \\ \hat{\mathbf{w}} &= (X^T X)^{-1} X^T \mathbf{y} \end{aligned}$$

Convert from vector summation form to matrix form

Left multiply by $(X^T X)^{-1}$

Cancel

$$\begin{aligned} X^T X &= [\mathbf{x}_1 \ \mathbf{x}_2 \ \cdots \ \mathbf{x}_n] \begin{bmatrix} \mathbf{x}_1^T \\ \mathbf{x}_2^T \\ \vdots \\ \mathbf{x}_n^T \end{bmatrix} \\ X^T X &= x_1 x_1^T + x_2 x_2^T + \cdots + x_n x_n^T \\ &= \sum_{i=1}^n x_i x_i^T \\ X^T \mathbf{y} &= [\mathbf{x}_1 \ \mathbf{x}_2 \ \cdots \ \mathbf{x}_n] \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} \\ X^T \mathbf{y} &= y_1 x_1 + y_2 x_2 + \cdots + y_n x_n \\ &= \sum_{i=1}^n x_i y_i \end{aligned}$$

Problem 4.2

Linear algebra facts for linear regression

We saw in lecture on Linear Regression that the closed form expression for linear regression without an offset involves the term $(X^T X)^{-1}$.

- (a) Is it true that the matrix $X^T X$ is always symmetric and positive semidefinite?

Yes. Symmetry can be checked by computing the transpose. For any vector u , we have $u^T X^T X u = \|Xu\|_2^2 \geq 0$.

- (b) State and prove the connection between the nullspace of X and the nullspace of $X^T X$. That is, your statement should look like one of the following: $\text{Null}(X) \subseteq \text{Null}(X^T X)$, or $\text{Null}(X) \supseteq \text{Null}(X^T X)$ or $\text{Null}(X) = \text{Null}(X^T X)$.

We have, $\text{Null}(X) = \text{Null}(X^T X)$. Let $v \in \text{Null}(X)$. Then, one can check that $X^T X v = 0$, leading to $v \in \text{Null}(X^T X)$, which proves $\text{Null}(X) \subseteq \text{Null}(X^T X)$. For the other direction, let $0 \neq v \in \text{Null}(X^T X)$. Then, $0 = v^T X^T X v = \|Xv\|_2^2$, which implies $v \in \text{Null}(X)$. Therefore, $\text{Null}(X^T X) \subseteq \text{Null}(X)$, which finishes the proof.

(c) Is it true that $X^\top X$ is always invertible?

No, this isn't always the case. Since $\text{Null}(X) = \text{Null}(X^\top X)$ (see the answer to the previous question), the matrix $X^\top X$ is not invertible if X has a non-empty nullspace.

(d) Based on the above fact about the connection between the nullspaces of X and $X^\top X$ and the expression for linear regression without an offset (that we referred to two problems above), justify the use of “tall skinny” data matrix X as opposed to a “short wide” matrix X .

If X is “short and wide”, it has a non-empty nullspace. Therefore, $X^\top X$ is not invertible.

(e) The columnspace and rowspace of $X^\top X$ are the same, and are equal to the rowspace of X . (Hint: Use the relationship between nullspace and rowspace.)

$X^\top X$ is symmetric, and from the previous parts, we have $\text{Row}(X^\top X) = \text{Col}((X^\top X)^\top) = \text{Col}(X^\top X)$. By previous parts again, we have: $\text{Row}(X^\top X) = \text{Null}(X^\top X)^\perp = \text{Null}(X)^\perp = \text{Row}(X)$.

Questions?