

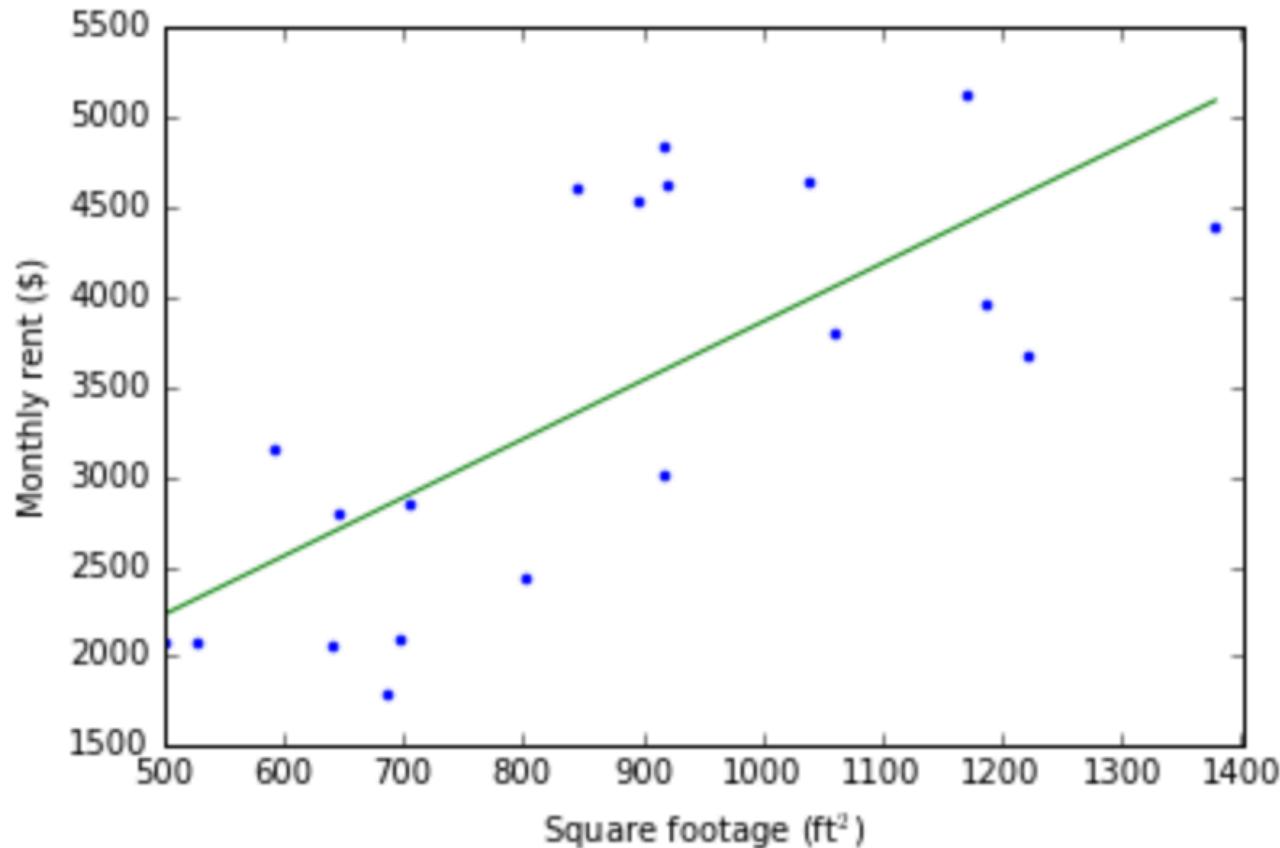


2018-01-17: Announcements

- HW #1 is uploaded to CCLE and due to Gradescope at 11:59pm on Monday, 22 Jan 2018.
- Please feel free to help answer questions on Piazza.
- Reminder that it is all okay to use Python 2.
- A/V system *ought* to be fine...
- Lecture notes continue to be updated and are labeled with the last date of update on CCLE.



2018-01-17: RECAP from 2018-01-10



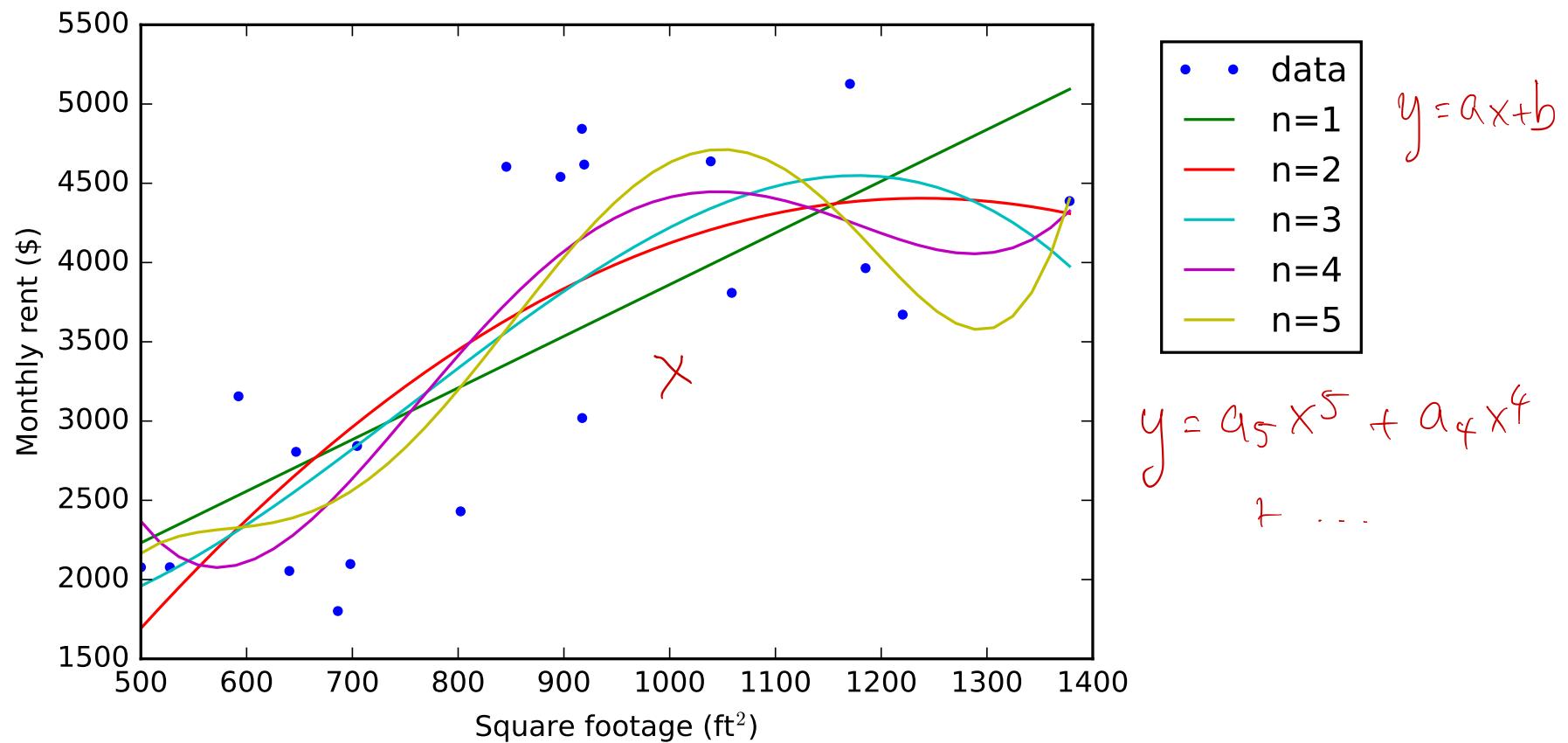
$$\mathcal{L} = \frac{1}{2} \sum_{i=1}^N (y^{(i)} - \theta^T \hat{\mathbf{x}}^{(i)})^2$$

$$y = \theta^T \hat{x}$$

$$\theta = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T Y$$

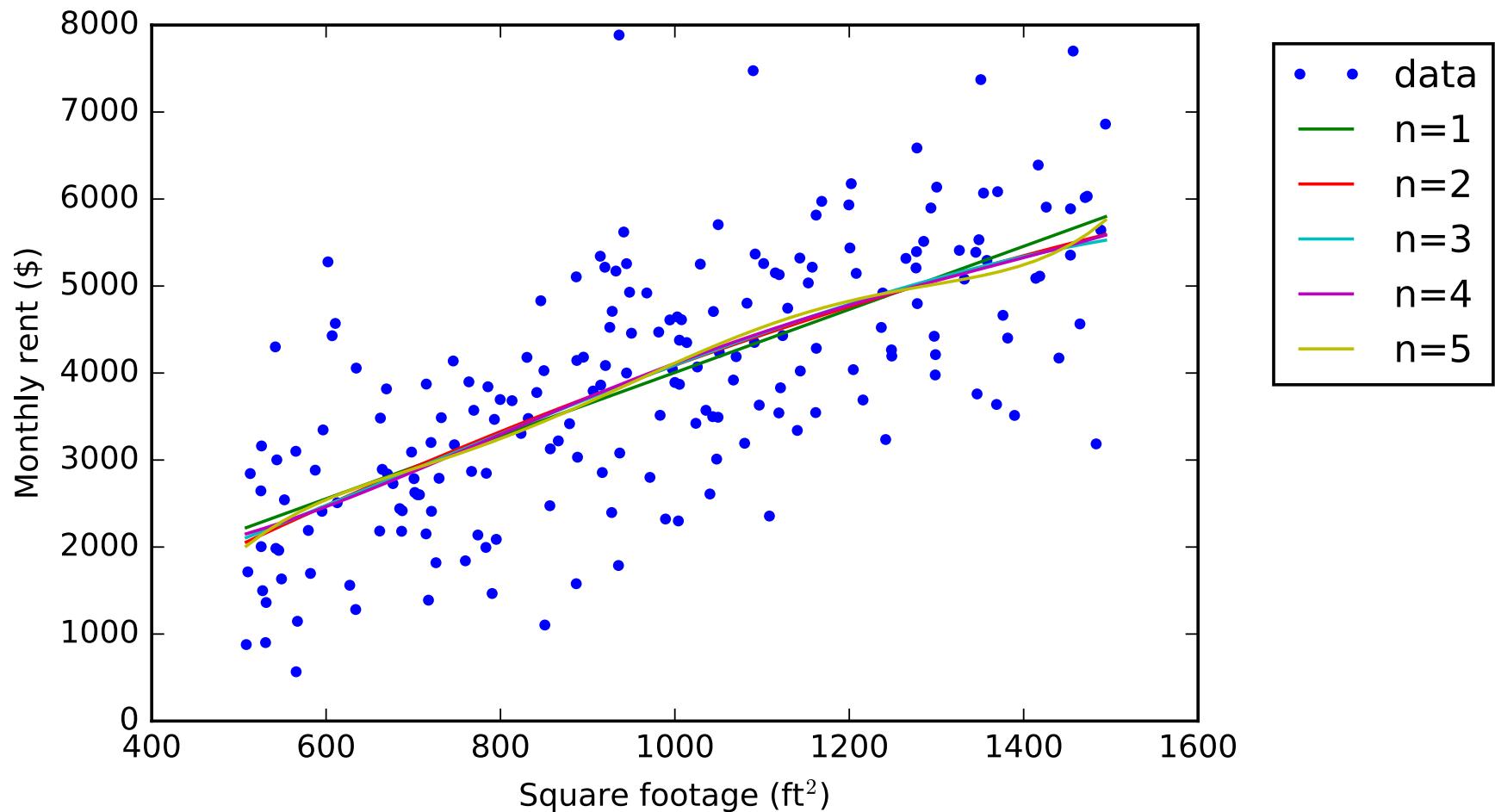


2018-01-17: RECAP from 2018-01-10



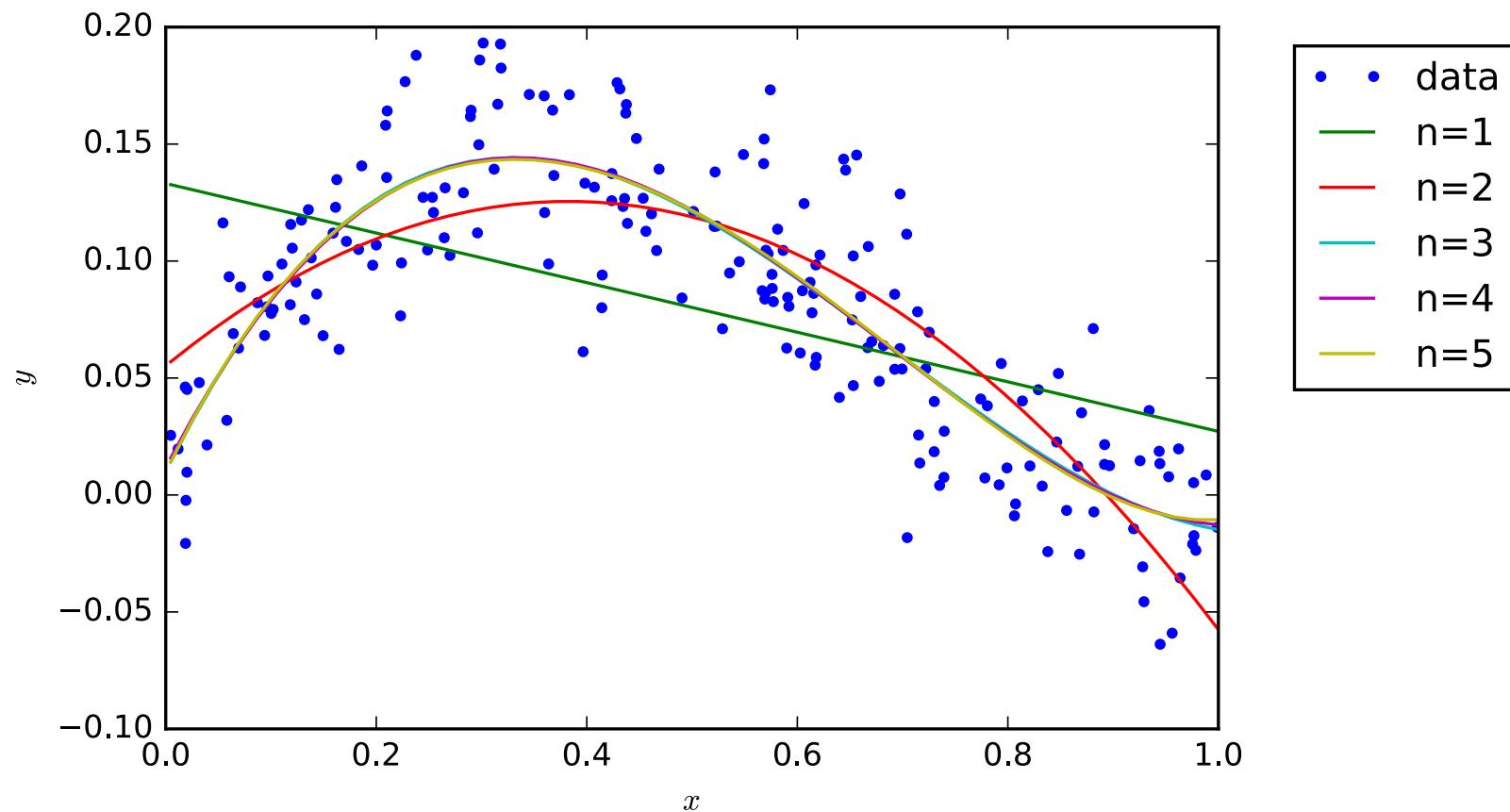


2018-01-17: RECAP from 2018-01-10





2018-01-17: RECAP from 2018-01-10

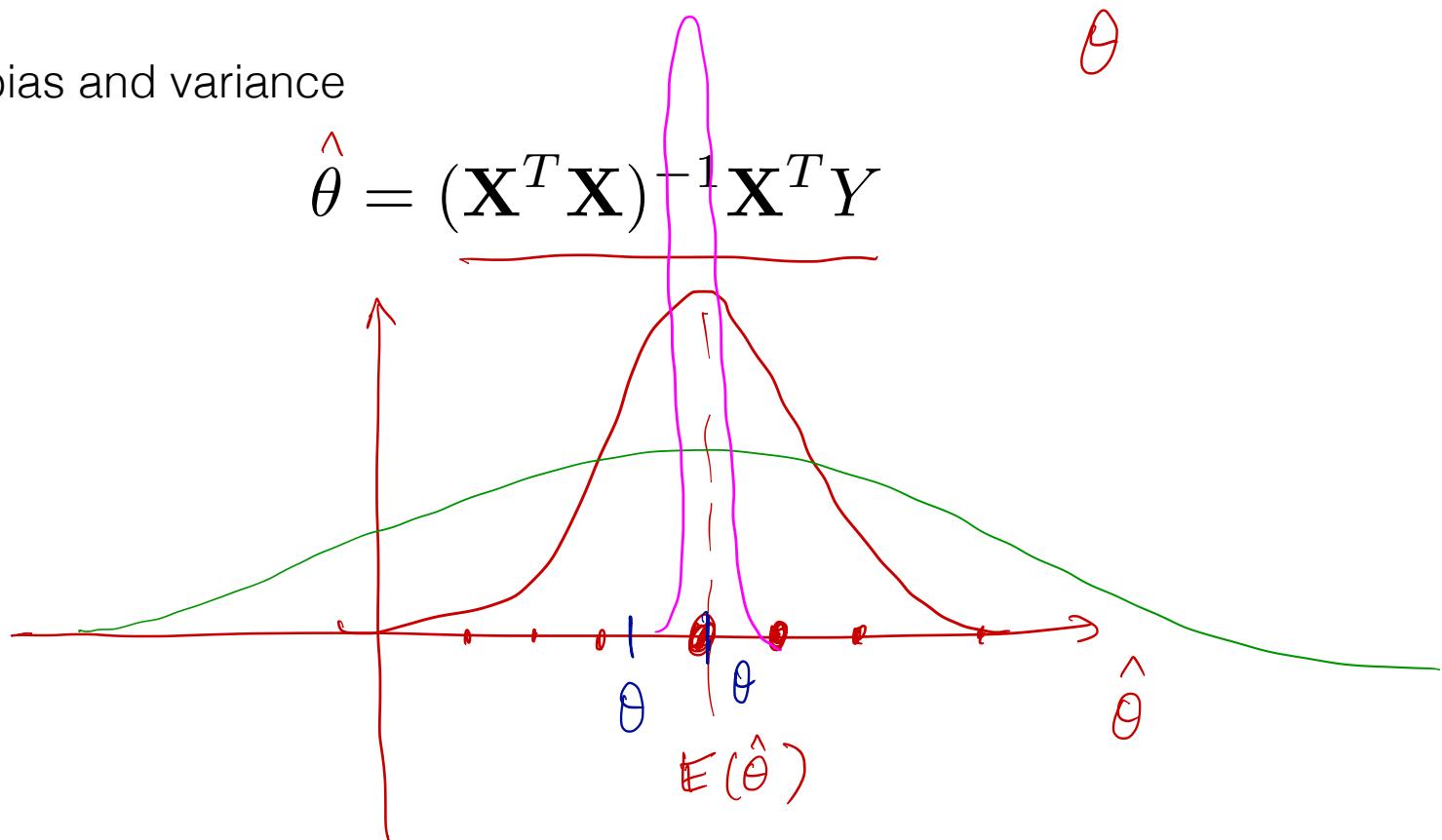




2018-01-17: RECAP from 2018-01-10

Estimator bias and variance

$$\hat{\theta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$$





Estimator bias

The *bias* of an estimator, $\hat{\theta}_m$, is defined as:

$$\text{bias}(\hat{\theta}_m) = \mathbb{E}(\hat{\theta}_m) - \theta$$

~~\mathbb{E}~~ —

Note:

- The expectation is taken over the data (where the randomness is introduced through samples $\mathbf{x}^{(i)}$).
- Informally, the bias measures how close $\hat{\theta}_m$ comes to estimating θ on average.
- An estimator is called *unbiased* if $\text{bias}(\hat{\theta}_m) = 0$.



Estimator bias example

Estimator bias (example)

$$x^{(i)} = \begin{cases} 0 & \text{w.p. } 1-\theta \\ 1 & \text{w.p. } \theta \end{cases}$$

- Let $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$ be iid samples from a Bernoulli distribution with mean θ .
- The *sample mean* estimator is given by:

$$\hat{\theta}_m = \frac{1}{m} \sum_{i=1}^m x^{(i)}$$

Intuitively, this ought to be a very reasonable estimator of the Bernoulli mean (or equivalently parameter).

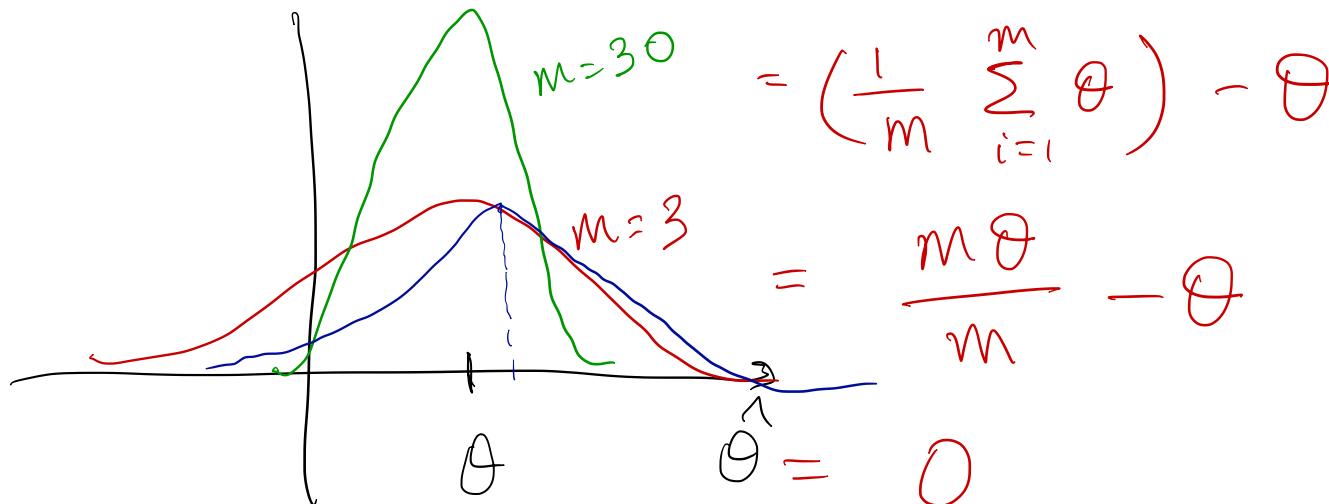


Estimator bias example

Then, the bias of this estimator is:

$$\begin{aligned}\text{bias}(\hat{\theta}_m) &= \mathbb{E}\left(\frac{1}{m} \sum_{i=1}^m x^{(i)}\right) - \theta \\ &= \left(\frac{1}{m} \sum_{i=1}^m \mathbb{E} x^{(i)}\right) - \theta\end{aligned}$$

$$\begin{aligned}\mathbb{E}(ax) &= a\mathbb{E}x \\ \mathbb{E}(x+y) &= \mathbb{E}x \\ &\quad + \mathbb{E}y\end{aligned}$$



$$= \left(\frac{1}{m} \sum_{i=1}^m \theta\right) - \theta$$

$$= \frac{m\theta}{m} - \theta$$

$$\theta = 0$$



Estimator variance

Estimator variance

An unbiased estimator may on average estimate θ , but any single instance of $\hat{\theta}$ may deviate from θ . The variance of $\hat{\theta}$ in predicting is called the variance of the estimator, denoted $\text{var}(\hat{\theta})$.

Some notes:

- This variability is due to the data samples that you get.
- The *standard error* of the estimator is $\sqrt{\text{var}(\hat{\theta})}$ and is commonly denoted $\text{SE}(\hat{\theta})$.



Estimator variance example

Estimator variance (example)

$$P(X, Y) = P(X) P(Y)$$

Let $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$ be iid samples from a Bernoulli distribution with mean θ . The *sample mean* estimator is given by:

$$\hat{\theta}_m = \frac{1}{m} \sum_{i=1}^m x^{(i)}$$

$$\text{var}(aX) = a^2 \text{var}(X)$$

$$\text{var}(X+Y) = \text{var}(X) + \text{var}(Y)$$

$$+ 2 \text{cov}(X, Y)$$

Then, the variance of this estimator is:

$$\begin{aligned}\text{var}(\hat{\theta}_m) &= \text{var}\left(\frac{1}{m} \sum_{i=1}^m x^{(i)}\right) \\ &= \frac{1}{m^2} \sum_{i=1}^m \text{var}(x^{(i)}) \\ &= \frac{1}{m^2} m (\theta(1-\theta)) \\ &= \frac{1}{m} \theta(1-\theta)\end{aligned}$$



Estimator variance example

Estimator variance (example)

A common metric used in evaluating the average of the data is the standard error of the mean. Given iid samples $x^{(1)}, \dots, x^{(m)}$, the standard error of the mean is given by:

$$\begin{aligned}\text{SE}(\hat{\mu}_m) &= \sqrt{\text{var}\left(\frac{1}{m} \sum_{i=1}^m x^{(i)}\right)} \\ &= \sqrt{\frac{1}{m^2} \sum_{i=1}^m \text{var}(x^{(i)})} \\ &= \sqrt{\frac{1}{m^2} m \text{var}(x^{(i)})} \\ &= \frac{\sigma}{\sqrt{m}}\end{aligned}$$

What does this mean intuitively?



Bias and mean-square error trade off

Trading off bias and variance

We may at times have to choose between a low-bias, high-variance estimator, or a high-bias, low-variance estimator.

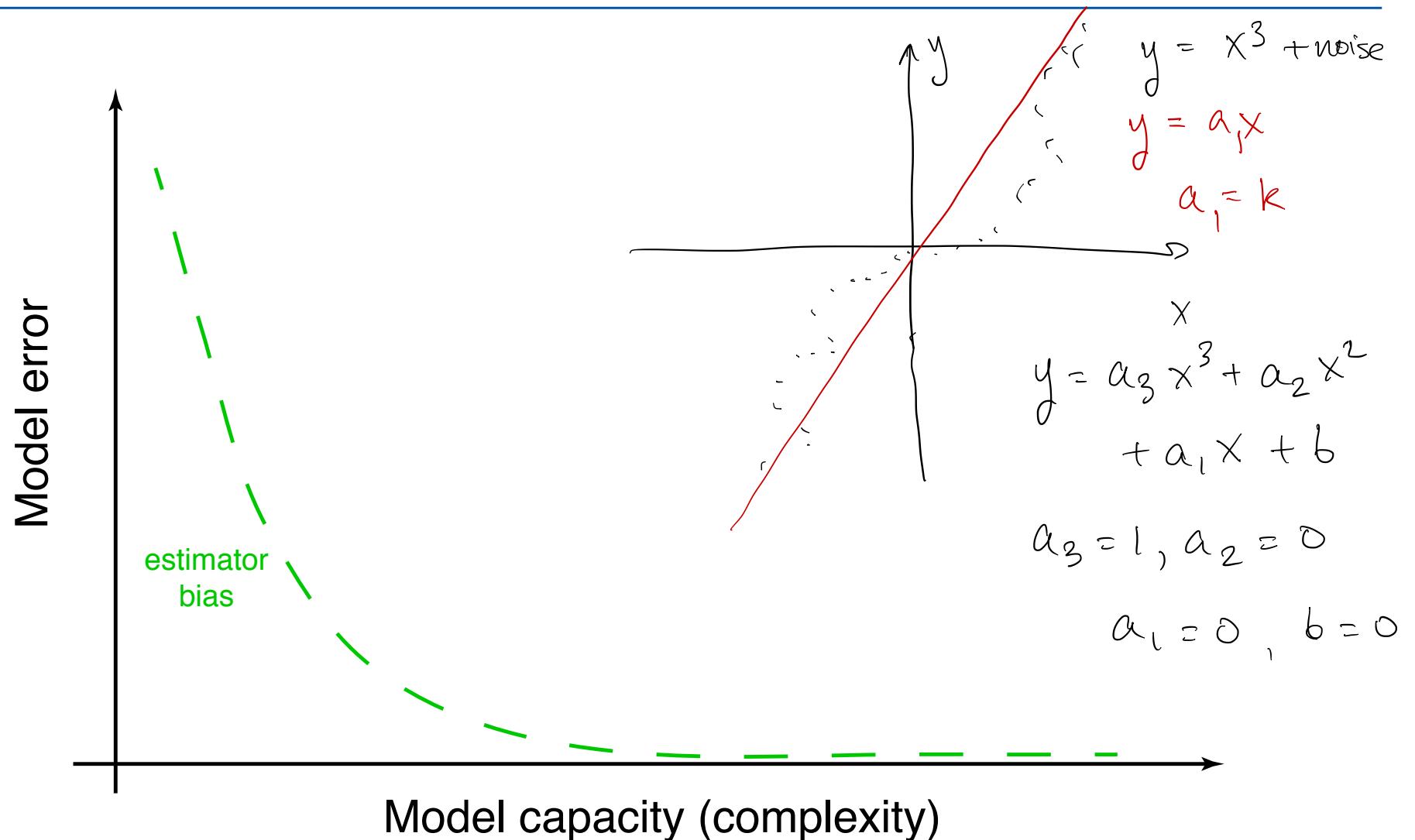
- The mean squared error (MSE) naturally trades these two off.

$$\begin{aligned}\text{MSE} &= \mathbb{E}(\hat{\theta}_m - \theta)^2 \\ &= \mathbb{E}(\hat{\theta}_m - \mathbb{E}\hat{\theta}_m + \mathbb{E}\hat{\theta}_m - \theta)^2 \\ &= \mathbb{E}(\hat{\theta}_m - \mathbb{E}\hat{\theta}_m)^2 + \mathbb{E} \left[2(\hat{\theta}_m - \mathbb{E}\hat{\theta}_m)(\mathbb{E}\hat{\theta}_m - \theta) \right] + \mathbb{E}(\mathbb{E}\hat{\theta}_m - \theta)^2 \\ &= \text{var}(\hat{\theta}_m) + 2\text{bias}(\hat{\theta}_m)\mathbb{E}(\hat{\theta}_m - \mathbb{E}\hat{\theta}_m) + \mathbb{E}\text{bias}(\hat{\theta}_m)^2 \\ &= \text{var}(\hat{\theta}_m) + \text{bias}(\hat{\theta}_m)^2\end{aligned}$$

- Thus, minimizing MSE will simultaneously minimize the bias and variance of the estimator.

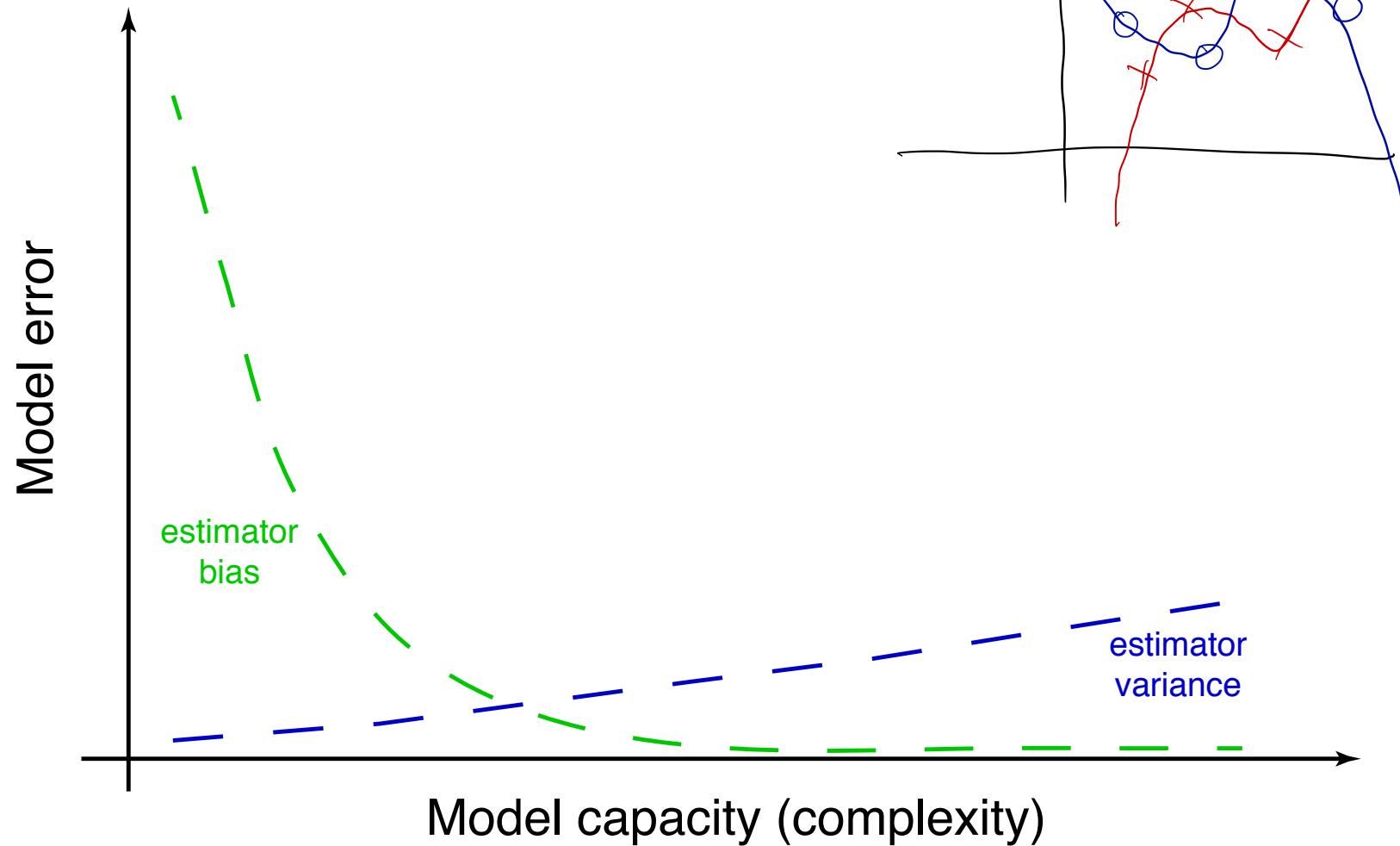


A picture to then have in mind



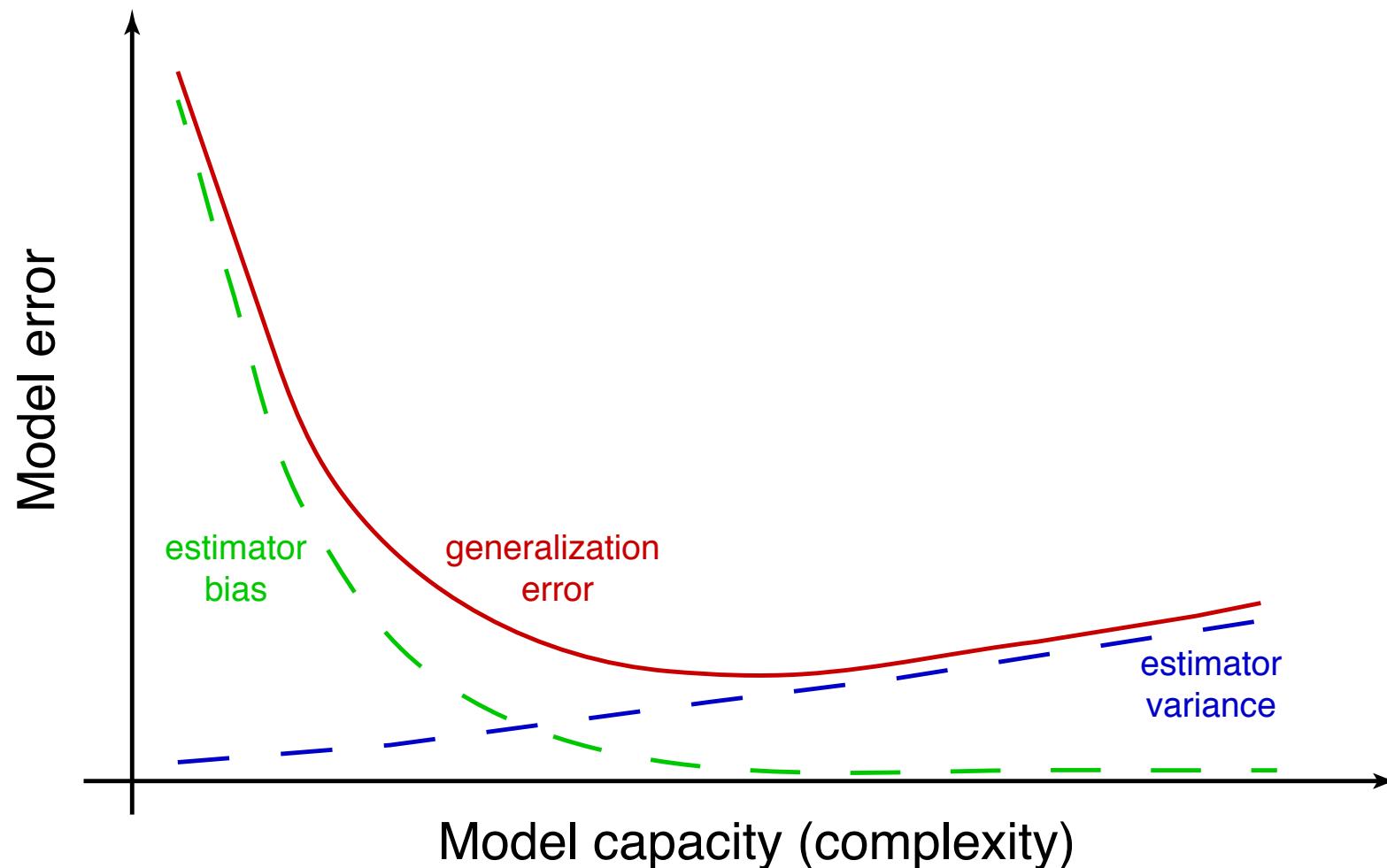


A picture to then have in mind



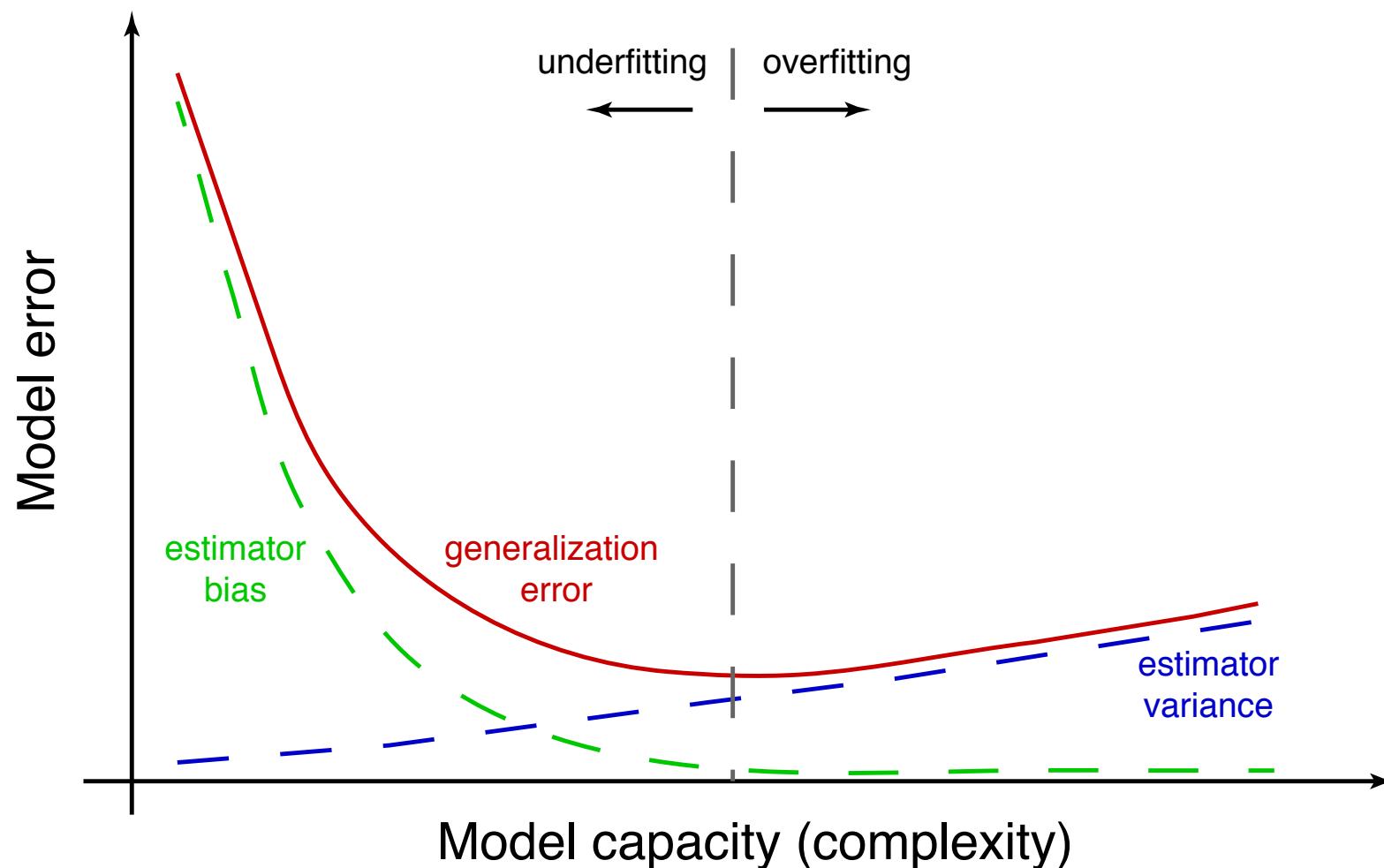


A picture to then have in mind





A picture to then have in mind

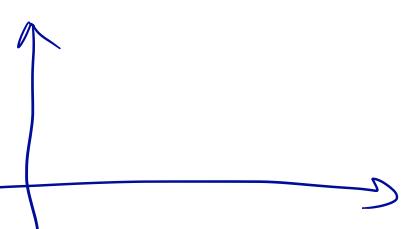




Picking a best model

From this picture, a reasonable way to pick a model is to assess its generalization error, and pick a setting of the parameters that results in minimal value.

A brief aside:

$$y = \alpha_1 x + b$$
$$y = \alpha_2 x^2 + \alpha_1 x + b$$


Sometimes there may be scenarios where dataset size is so limited that it is not possible to have a reasonably good test dataset. In these cases, we may utilize *model selection*, which is itself an entire field.

At a high level, the principles of model selection is to find a way to *penalize* the model for being overly complex. There is an art to designing some penalty terms.

Criterion for selecting models include:

- Bayes information criterion
- Akaike information criterion
- Deviance information criterion

My recommendation: if you have enough data to make a test set, do it.



Evaluating generalization error

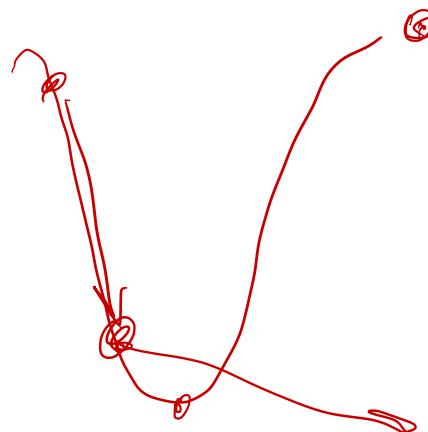
$$y = \theta^T X$$

Training, validation, and testing data

The standard for training, evaluating, and choosing models is to use different datasets for each step.

- **Training data** is data that is used to learn the parameters of your model.
- **Validation data** is data that is used to optimize the hyperparameters of your model. This avoids the potential of overfitting to nuances in the testing dataset.
- **Testing data** is data that is used to score your model.

Note, in many cases, testing and validation datasets are used interchangeably as datasets that are used to evaluate the model. In this scenario, hyperparameters would also be optimized using training data.





Evaluating generalization error

***k*-fold cross validation**

In a common scenario, you will be given a training dataset and a testing dataset. To train a model using this dataset, one common approach is k -fold cross validation. Then the procedure looks as follows:

- Let the training dataset contain N examples. 1000 $K=5$
- Then, split the data into k equal sets, each of N/k examples. Each of these sets is called a “fold.” 200 800 200
- $k - 1$ of the folds are datasets that are used to train the model parameters.
- The remaining fold is a testing dataset used to evaluate the model.
- You may repeatedly train the model by choosing which folds comprise the training folds and the testing fold.



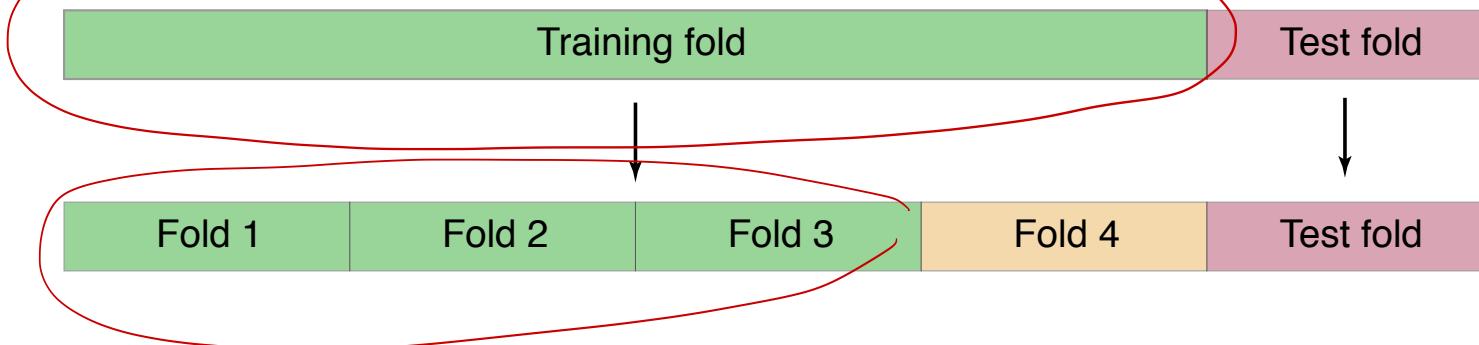
Evaluating generalization error



k-fold cross-validation with no hyperparameters



k-fold cross-validation with hyperparameters





k-fold cross validation from housing example

Routines for training models and testing data.

```
def train_models(x, y):
    # Fit a polynomial model up to N=5
    N = 5
    thetas = []
    for i in np.arange(N):
        xhat = np.vstack((x, np.ones_like(x))) if i == 0 else np.vstack((x**(i+1), xhat))
        thetas.append(np.linalg.inv(xhat.dot(xhat.T)).dot(xhat.dot(y)))

    return thetas

def test_models(x, y, thetas):
    N = 5
    errors = []
    for i in np.arange(N):
        xhat = np.vstack((x, np.ones_like(x))) if i == 0 else np.vstack((x**(i+1), xhat))
        errors.append(np.sqrt(np.sum((y - thetas[i].dot(xhat))**2) / len(y)))
    return errors
```



k-fold cross validation from housing example

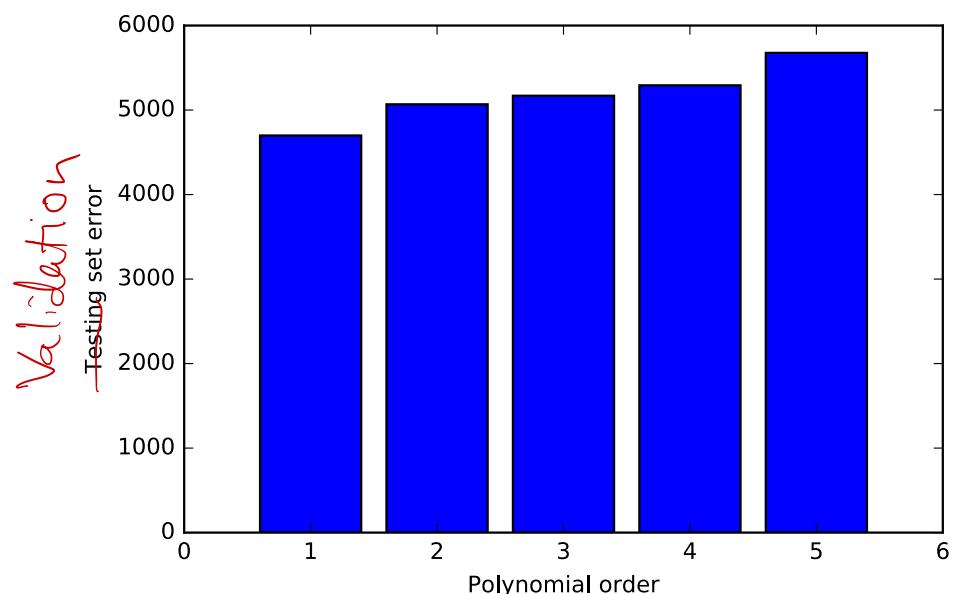
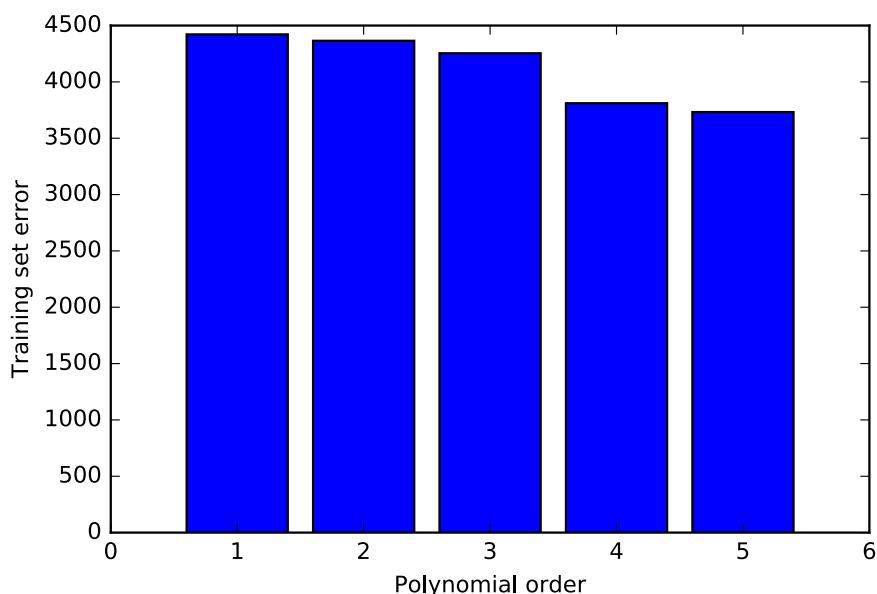
```
# (x,y) are the training data pairs from the supervised housing rent example
# num_train = len(x) is the number of examples
num_folds = 5
cv_idx = np.arange(num_train)          1000          200
np.random.shuffle(cv_idx)
fold_size = num_train // 5
train_error, test_error = (np.zeros(5), np.zeros(5))
for i in np.arange(num_folds):
    test_idx = cv_idx[i*fold_size:(i+1)*fold_size]
    train_idx = np.concatenate((cv_idx[:i*fold_size], cv_idx[(i+1)*fold_size:]))
    x_test = x[test_idx]
    y_test = y[test_idx]
    x_train = x[train_idx]
    y_train = y[train_idx]
    thetas = train_models(x_train, y_train)
    train_error += test_models(x_train, y_train, thetas)
    test_error += test_models(x_test, y_test, thetas)

print train_error
print test_error
f = plt.figure()
ax = f.gca()
ax.bar(np.arange(5) + 1, train_error, width=0.8, align='center')
ax.set_xlabel('Polynomial order')
ax.set_ylabel('Training set error')
f.savefig('ml-basics_cv-train.pdf')
f = plt.figure()
ax = f.gca()
ax.bar(np.arange(5) + 1, test_error, width=0.8, align='center')
ax.set_xlabel('Polynomial order')
ax.set_ylabel('Testing set error')
f.savefig('ml-basics_cv-test.pdf')
```

```
[ 4421.38242236  4364.65797472  4254.01787959  3810.5847096  3731.95920537]
[ 4698.83840383  5067.8036503   5169.2079499   5292.40392825  5676.26492547]
```



k-fold cross validation from housing example



*Validation
Testing set error*



Other types of optimization

We've talked about examples where we want to *minimize* a mean-square error or distance metric.

Another metric that we may want to minimize is the *probability of having observed the data*. In this framework, the data is modeled to have some distribution with parameters. We choose the parameters to maximize the probability of having observed our training data.

$$x^{(i)} = \begin{cases} 0 & \text{w.p. } 1 - \theta \\ 1 & \text{w.p. } \theta \end{cases}$$

HTHHHTTHT

$$\text{Model 1: } \theta = 1$$

$$\text{Model 2: } \theta = 0.5$$

$$\text{Model 3: } \theta = 0.75$$

$$\left\{ \begin{array}{l} \text{Model 1: } 1^4 \cdot 0^4 = 0 \\ \text{Model 2: } (0.5)^8 = 0.0039 \\ \text{Model 3: } (0.75)^4 (0.25)^4 = 0.00124 \end{array} \right.$$



Maximum-likelihood introduction

H H H H H H

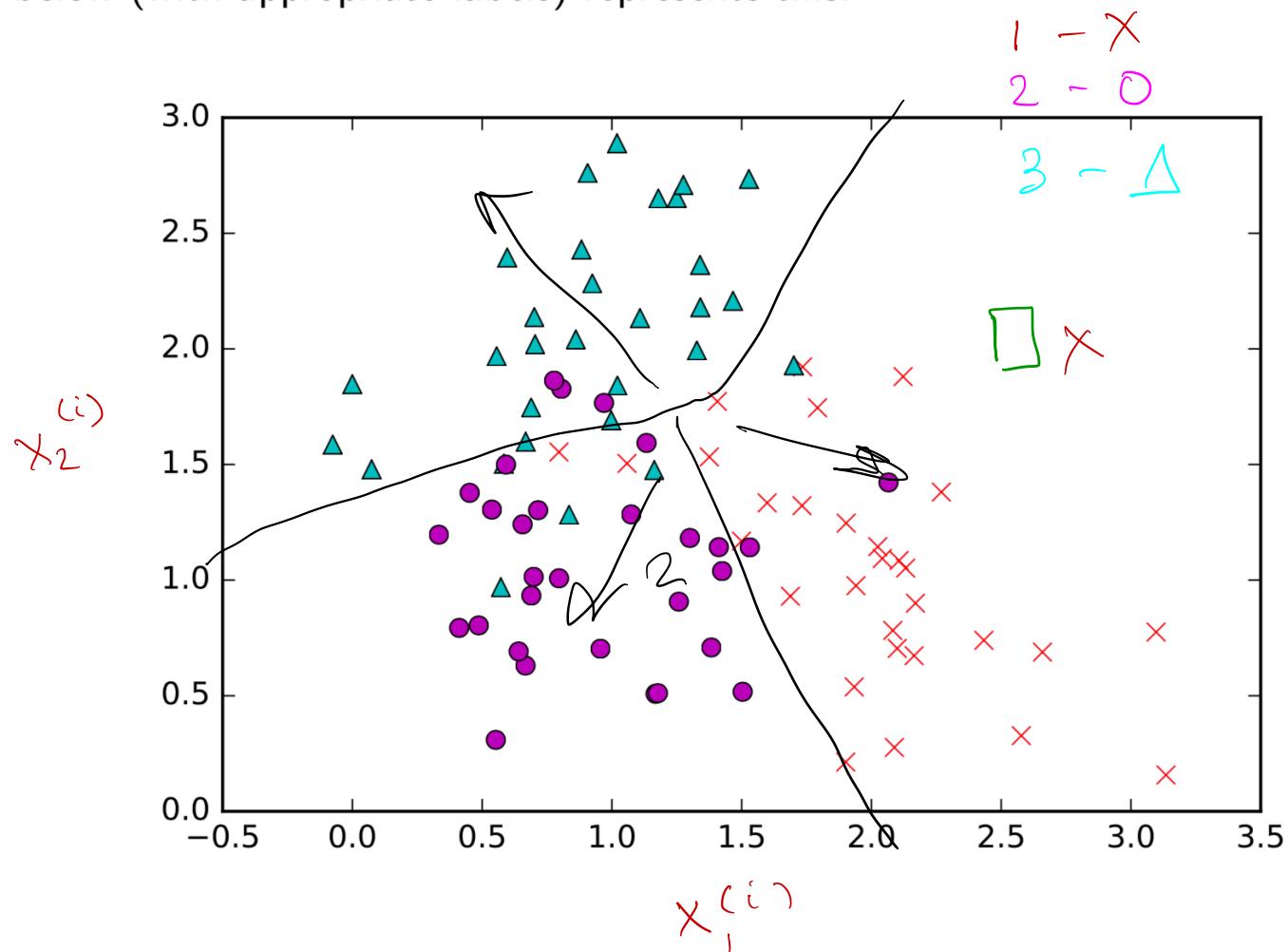
Pr() under Model 1 : $1^5 \cdot (0)^0 = 1$



Maximum-likelihood estimation

Example of ML estimation

Say we receive paired data $\{\mathbf{x}^{(i)}, y^{(i)}\}$ where $\mathbf{x}^{(i)} \in \mathbb{R}^2$ is a data point that belongs to one of three classes, $y^{(i)} \in \{1, 2, 3\}$. Classes 1, 2, 3 denote the three possible classes (or labels) that a data point could belong to. The drawing below (with appropriate labels) represents this:





What is our goal in modeling?



What is our goal in modeling?



Chain rule for probability

$$\Pr(A = a) = p_A(a) = p(a)$$

$$\Pr(B = b) = p_B(b) = p(b)$$

$$\Pr\{A = a \text{ and } B = b\} = p_{A,B}(a, b) = p(a, b)$$

$$\begin{aligned}\Pr\{A = a \text{ and } B = b\} &= \Pr\{A = a\} \Pr\{B = b \text{ given } A = a\} \\ &= p(a) p(b | a) \\ &= p(b) p(a | b)\end{aligned}$$

$$\begin{aligned}p(a, b, c) &= p(c) p(b | c) p(a | b, c) \\ &= p(b) p(a | b) p(c | a, b) \\ &= p(a) p(c | a) p(b | a, c)\end{aligned}$$



Chain rule for probability

$$p(b, c | d, e) = \frac{p(b, c, d, e)}{p(d) p(e | d)}$$

$p(e, d)$

A blue curly brace underlines the terms $p(d)$ and $p(e | d)$.

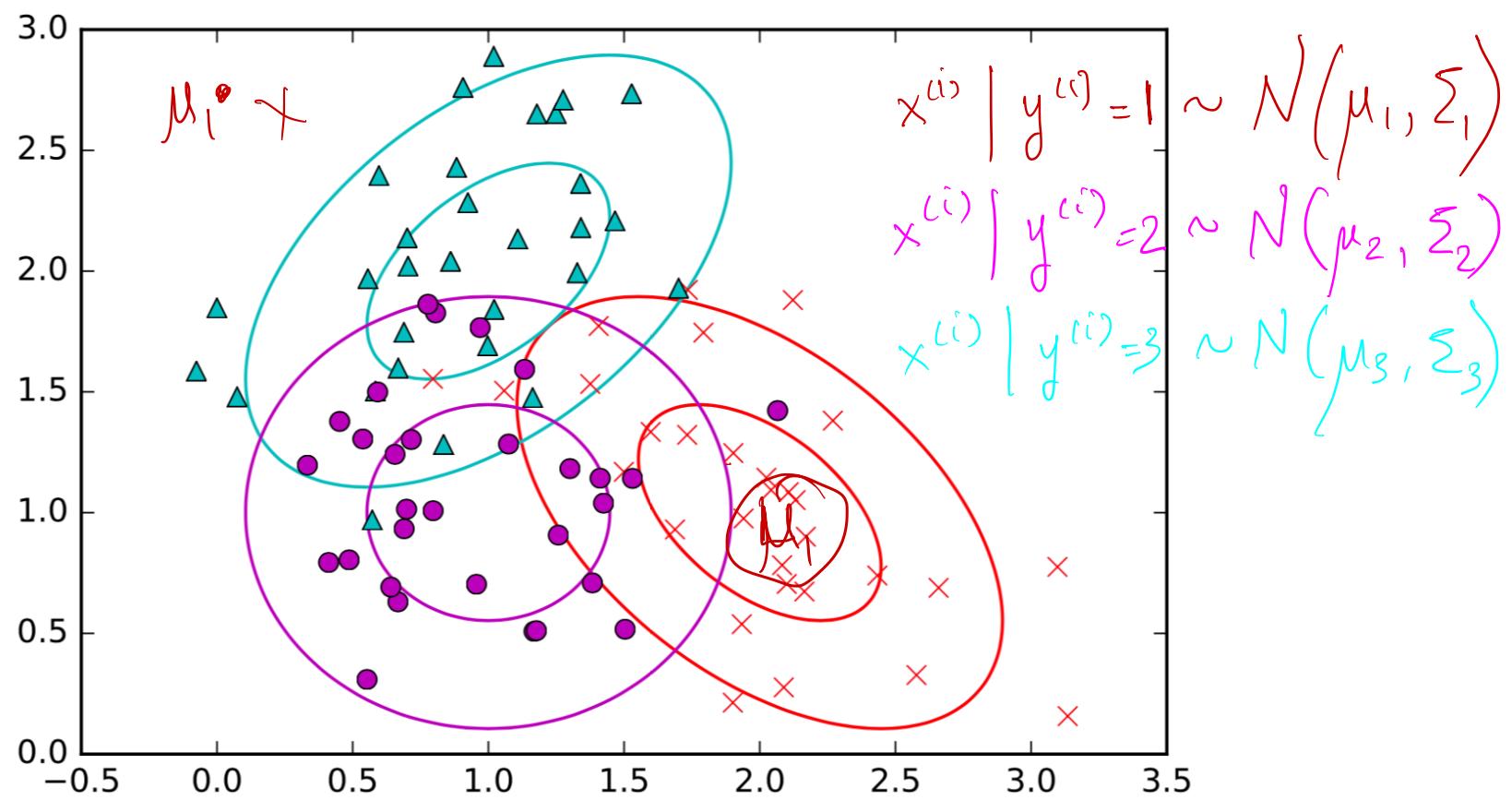
$$\begin{aligned} p(b, c, d, e) &= p(e, d) p(b, c | d, e) \\ &= p(e, d) p(b | d, e) \\ &\quad p(c | b, d, e) \end{aligned}$$



Maximum-likelihood estimation

Example of ML estimation (cont.)

We may want to model each cluster as being *generated* by a multivariate Gaussian distribution. Concretely, in this example, there are three Gaussian clusters, each one describing the distribution of points for that class.





Maximum-likelihood estimation

Example of ML estimation (cont.)

The model setup is as follows:

- Each data point $\mathbf{x}^{(i)}$ belongs to class $y^{(i)}$.
- Each class y_i is parametrized according to a distribution.

$$\begin{aligned}\mathbf{x}^{(i)} | y^{(i)} = 1 &\sim \mathcal{N}(\mu_1, \Sigma_1) \\ \mathbf{x}^{(i)} | y^{(i)} = 2 &\sim \mathcal{N}(\mu_2, \Sigma_2) \\ \mathbf{x}^{(i)} | y^{(i)} = 3 &\sim \mathcal{N}(\mu_3, \Sigma_3)\end{aligned}$$

Thus, the parameters we can choose to optimize our model are

$$\theta = \{\mu_1, \Sigma_1, \mu_2, \Sigma_2, \mu_3, \Sigma_3\}.$$

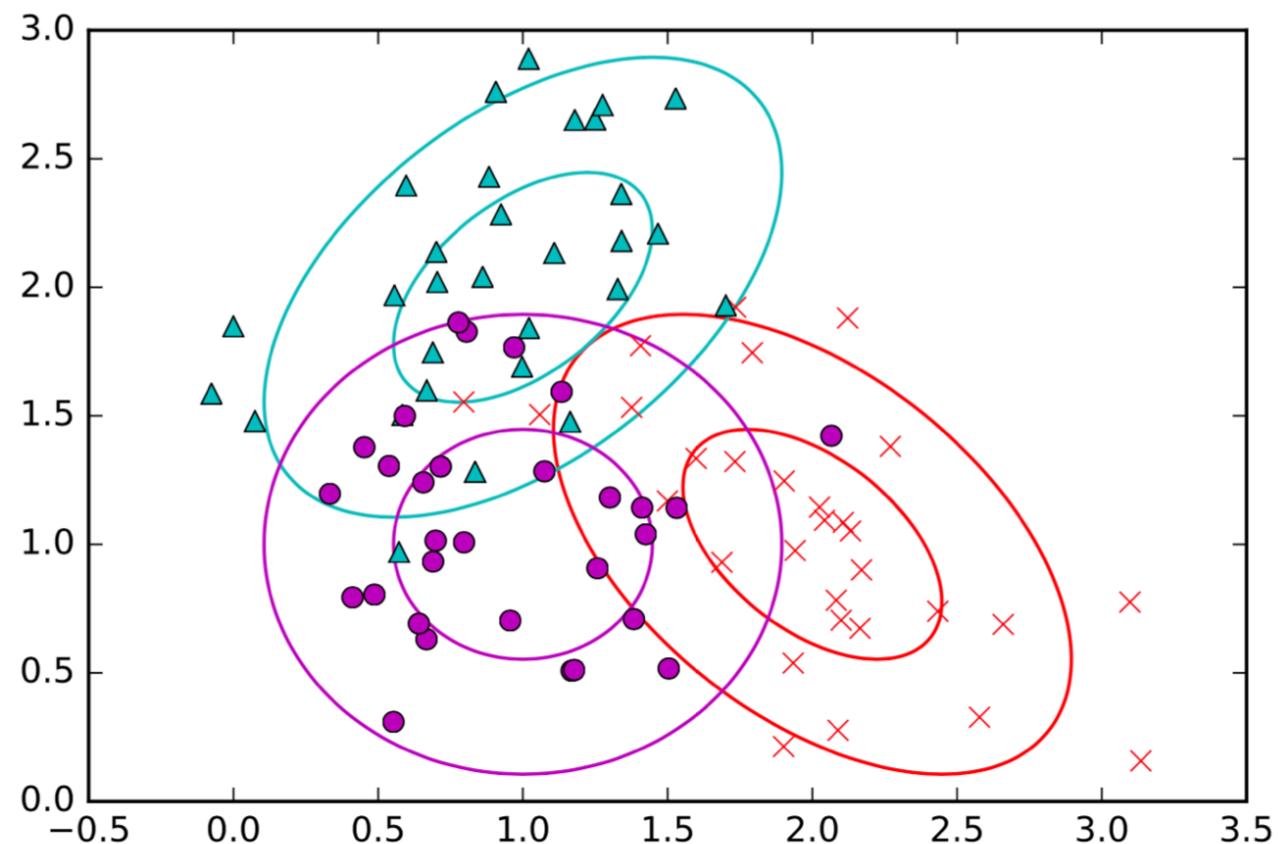
$$\Pr(y^{(i)} = j) = \frac{1}{3}$$

- We'll assume all classes are equally probable a priori (so that maximum likelihood estimation and maximum a posteriori estimation are equivalent).
- Finally, we'll assume each data point is independent, so that we can easily write out probabilities, i.e.,

$$p\left(\{\mathbf{x}^{(i)}, y^{(i)}\}, \{\mathbf{x}^{(j)}, y^{(j)}\}\right) = p\left(\mathbf{x}^{(i)}, y^{(i)}\right) p\left(\mathbf{x}^{(j)}, y^{(j)}\right)$$



Maximum-likelihood estimation





Maximum-likelihood estimation

Example of ML estimation (cont.)

We'd like to maximize the likelihood of having seen the dataset. This can be written as

$$\mathcal{L} = p\left(\{\mathbf{x}^{(1)}, y^{(1)}\}, \{\mathbf{x}^{(2)}, y^{(2)}\}, \dots, \{\mathbf{x}^{(N)}, y^{(N)}\}\right)$$

$$= \prod_{i=1}^N p(x^{(i)}, y^{(i)})$$

$$= \prod_{i=1}^N p(x^{(i)} | y^{(i)}) p(y^{(i)})$$

$$\underline{\log \mathcal{L}} = \sum_{i=1}^N \left[\log p(y^{(i)}) + \log p(x^{(i)} | y^{(i)}) \right]$$

$$= k + \sum_{i=1}^N \log p(x^{(i)} | y^{(i)})$$



Maximum-likelihood estimation



Maximum-likelihood estimation

Example of ML estimation (cont.)

Now, to simplify the probabilities, we evaluate $\log p(\mathbf{x}_i | y_i)$. We first observe:

$$\begin{aligned}\underbrace{\log p(\mathbf{x}^{(i)} | y^{(i)} = j)} &= \log \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma_j|^{\frac{1}{2}}} \exp \left(-\frac{1}{2} (\mathbf{x}^{(i)} - \mu_j)^T \Sigma_j^{-1} (\mathbf{x}^{(i)} - \mu_j) \right) \\ &= (2\pi)^{-n/2} |\Sigma_j|^{-1/2} \exp \left(\downarrow \right) \\ &= -\frac{n}{2} \log 2\pi - \frac{1}{2} \log |\Sigma_j| - \frac{1}{2} (\mathbf{x}^{(i)} - \mu_j)^T \Sigma_j^{-1} (\mathbf{x}^{(i)} - \mu_j) \\ &= k_2 - \frac{1}{2} \log |\Sigma_j| - \frac{1}{2} (\mathbf{x}^{(i)} - \mu_j)^T \Sigma_j^{-1} (\mathbf{x}^{(i)} - \mu_j)\end{aligned}$$



Maximum-likelihood estimation

$$\log \mathcal{L} = k_3 + \sum_{i=1}^N -\frac{1}{2} \log |\Sigma_{y^{(i)}}| - \frac{1}{2} (x^{(i)} - \mu_{y^{(i)}})^\top \Sigma_{y^{(i)}}^{-1} (x^{(i)} - \mu_{y^{(i)}})$$

$$k_3 = N k_2 + N \log \frac{1}{3}$$



Maximum-likelihood estimation

$$\frac{\partial \log \mathcal{L}}{\partial \mu_j} = \frac{\partial}{\partial \mu_j} \left[\sum_{i:y^{(i)}=j} -\frac{1}{2} (\mathbf{x}^{(i)} - \mu_j)^T \Sigma_j^{-1} (\mathbf{x}^{(i)} - \mu_j) \right]$$

$$= \frac{-1}{2} \frac{\partial}{\partial \mu_j} \left[\sum_{i:y^{(i)}=j} \mathbf{x}^{(i)\top} \Sigma_j^{-1} \mathbf{x}^{(i)} - \mathbf{x}^{(i)\top} \Sigma_j^{-1} \mu_j \right. \\ \left. - \mu_j^\top \Sigma_j^{-1} \mathbf{x}^{(i)} + \mu_j \Sigma_j^{-1} \mu_j \right]$$

$$\frac{\partial}{\partial \mu_j} \mu_j^\top \mathbf{a} = \mathbf{a}$$

$$\frac{\partial}{\partial x} x^\top A x \\ = (A + A^\top)x \\ = -\frac{1}{2} \frac{\partial}{\partial \mu_j} \left[\sum_{i:y^{(i)}=j} \mathbf{x}^{(i)\top} \Sigma_j^{-1} \mathbf{x}^{(i)} - 2\mu_j^\top \Sigma_j^{-1} \mathbf{x}^{(i)} \right. \\ \left. + \mu_j \Sigma_j^{-1} \mu_j \right]$$

$$= -\frac{1}{2} \sum \left[0 - 2 \Sigma_j^{-1} \mathbf{x}^{(i)} + 2 \Sigma_j^{-1} \mu_j \right] \\ = 0$$



Maximum-likelihood estimation

$$\sum_{i:y^{(i)}=j} \left(\sum_j^{-1} x^{(i)} \right) = \sum_{i:y^{(i)}=j} \left(\sum_j^{-1} \mu_j \right)$$

$$\sum_{i:y^{(i)}=j} x^{(i)} = N_j \mu_j$$

↙

$$N_j = \sum_{i:y^{(i)}=j} 1$$

$$\mu_j = \frac{1}{N_j} \sum_{i:y^{(i)}=j} x^{(i)}$$



Maximum-likelihood estimation

Example of ML estimation (cont.)

We next take the derivative with respect to Σ , using the following facts (for more information, see “Tools” lecture notes on derivatives w.r.t. matrices):

$$\begin{aligned}\frac{\partial}{\partial \Sigma} \text{tr}(\Sigma^{-1} \mathbf{A}) &= -\Sigma^{-T} \mathbf{A}^T \Sigma^{-T} \\ \frac{\partial}{\partial \Sigma} \log |\Sigma| &= \Sigma^{-T}\end{aligned}$$

Now, differentiating:

$$\frac{\partial \log \mathcal{L}}{\partial \Sigma_j} = \frac{\partial}{\partial \Sigma_j} \left[\sum_{i:y^{(i)}=j} -\frac{1}{2} \log |\Sigma_j| - \frac{1}{2} (\mathbf{x}^{(i)} - \mu_j)^T \Sigma_j^{-1} (\mathbf{x}^{(i)} - \mu_j) \right]$$



Maximum-likelihood estimation



Maximum-likelihood estimation

Example of ML estimation (cont.)

Solving for the optimal parameters μ_j and Σ_j for all classes, we arrive at the following solution:

```
# x and y are the data and their corresponding labels
# x is 2 x num_examples of the data points
# y is 1 x num_examples of the data labels

class_idxs = [np.where(y == j)[0] for j in np.unique(y)]
means = []
covs = []

for this_class in class_idxs:
    means.append(np.mean(x[:, this_class], axis=1))
    covs.append(np.cov(x[:, this_class]))
```



Maximum-likelihood estimation

Example of ML estimation (cont.)

Solving for the optimal parameters μ_j and Σ_j for all classes, we arrive at the following solution:

```
f = plt.figure()
ax = f.gca()
ax.plot(data1[0,:], data1[1,:], marker='x', linestyle=' ', color=class_colors[0])
ax.plot(data2[0,:], data2[1,:], marker='^', linestyle=' ', color=class_colors[1])
ax.plot(data3[0,:], data3[1,:], marker='o', linestyle=' ', color=class_colors[2])

for i in np.arange(len(class_idxs)):
    vals, vecs = eigsorted(covs[i])
    theta = np.degrees(np.arctan2(*vecs[:,0][::-1]))
    for nstd in np.arange(3):
        w, h = 2 * nstd * np.sqrt(vals)
        e = Ellipse(xy=(means[i][0], means[i][1]),
                    width=w, height=h,
                    angle=theta, color=class_colors[i])
        e.set_facecolor('none')
        ax.add_artist(e)

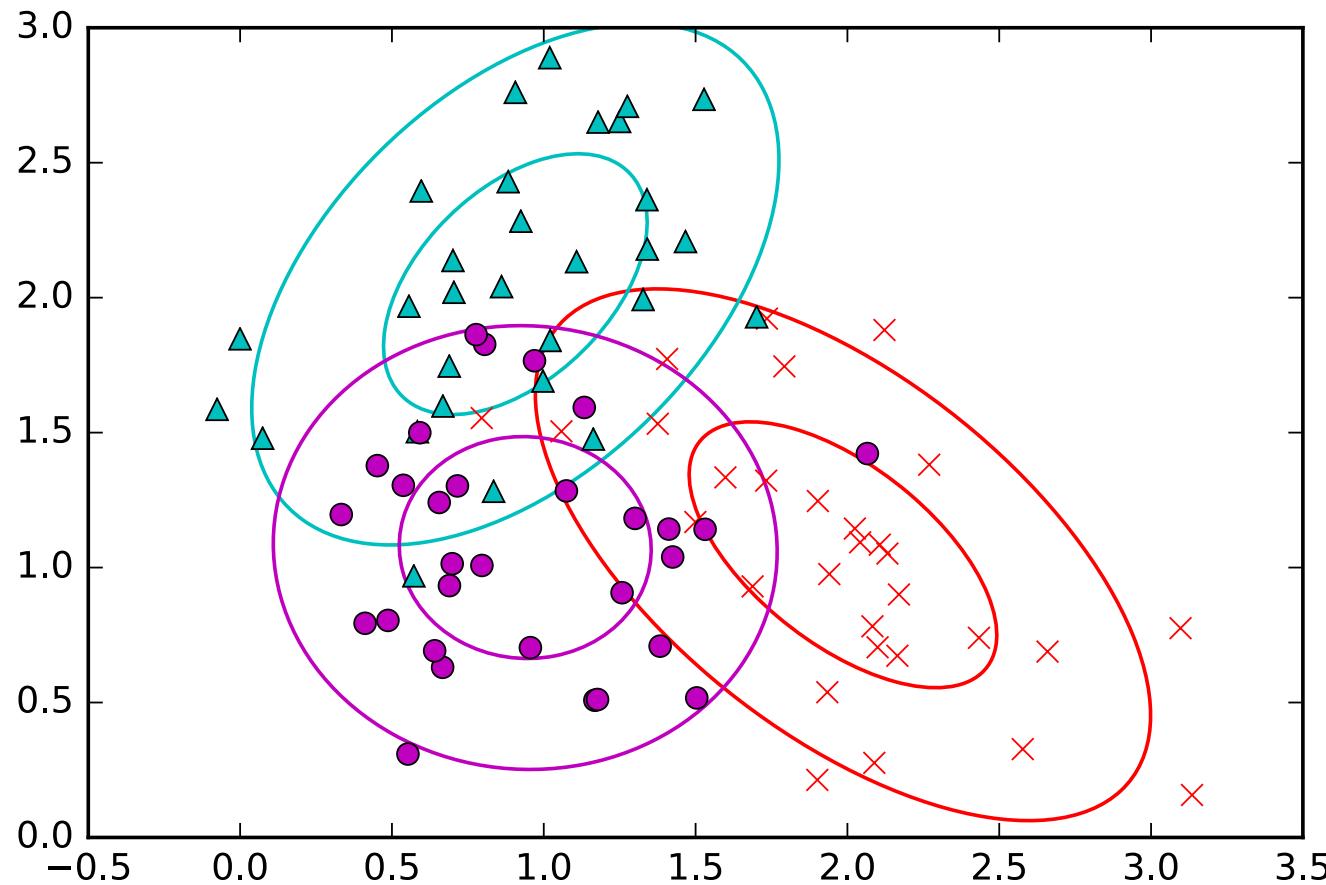
f.savefig('ml-basics_maxlik-data_fit-ellipsoids.pdf')
```



Maximum-likelihood estimation

Example of ML estimation (cont.)

Solving for the optimal parameters μ_j and Σ_j for all classes, we arrive at the following solution:





Maximum-likelihood estimation

Example of ML estimation (cont.)

Imagine now a new point \mathbf{x} comes in that we'd like to classify as being in one of three classes. To do so, we'd like to calculate: $\Pr(y = j|\mathbf{x})$ and pick the j that maximizes this probability. We'll denote this probability $p(j|\mathbf{x})$. We'll also assume, as earlier, that the classes are equally probable, i.e., $\Pr(y = j)$ is the same for all j .



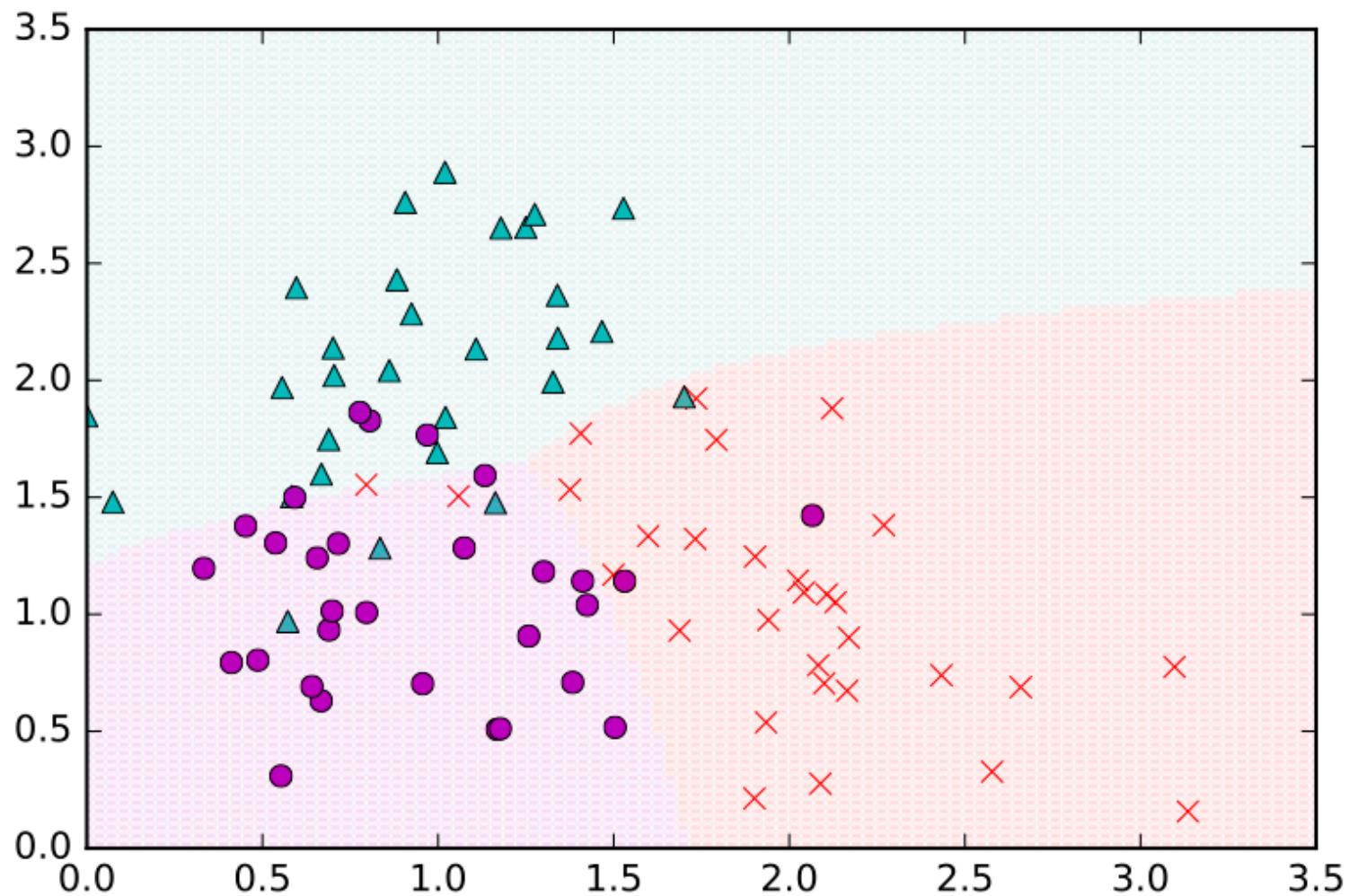
Maximum-likelihood estimation

```
f = plt.figure()
ax = f.gca()
ax.plot(data1[0,:], data1[1,:], marker='x', linestyle='', color=class_colors[0])
ax.plot(data2[0,:], data2[1,:], marker='^', linestyle='', color=class_colors[1])
ax.plot(data3[0,:], data3[1,:], marker='o', linestyle='', color=class_colors[2])

for (i,xx) in enumerate(np.linspace(0, 3.5, 100)):
    for yy in np.linspace(0, 3.5, 100):
        xn = np.array((xx,yy))
        pmax = -np.Inf
        decoded_class = np.nan
        for c in np.arange(len(class_idxs)):
            class_prob = -np.log(np.linalg.det(covs[c])) - (xn-means[c]).T.dot(np.linalg.inv(covs[c])).dot(xn-means[c])
            if class_prob > pmax:
                pmax = class_prob
                decoded_class = c
        ax.plot(xx, yy, marker='.', color=class_colors[decoded_class], alpha=0.05)
ax.set_xlim(0,3.5)
f.savefig('ml-basics_maxlik-data_fit-ellipsoids_colorized.pdf')
```



Maximum-likelihood estimation





Where we go from here

Even more cost functions

There are even more cost functions that we could use. We'll encounter some others in this class. Others include:

- MAP estimation.
- KL divergence.
- Maximize an approximation of the distribution.

In machine learning, it is important to arrive at an appropriate model and cost function. After that, it's important to know how to optimize it. In our examples here, the models were simple enough that we could differentiate and set the derivative equal to zero. In future lectures, we'll discuss more general ways to learn parameters of models when they are not so simple.



EXTRA SLIDES
