



Lecture 2: Reviewing the basics of machine learning

Announcements:

- We sent out an Announcement on Bruin Learn on how to sign up for Gradescope and Piazza, as well as a poll on discussion section scheduling.
- We are still setting the discussion sections and OH based on the polls. We will send a Bruin Learn announcement when these are set (likely this evening; tomorrow at the latest).
- We uploaded a few supporting files to Bruin Learn: a pdf on setting up Python and Jupyter Notebooks, and a pdf on using Google Colab (optional). You are welcome to use whatever setup you desire. I would personally install anaconda and create conda virtual environments.
- We also uploaded the formal notes of the class and exams dating back to 2018.
- Any questions on the syllabus or logistics?

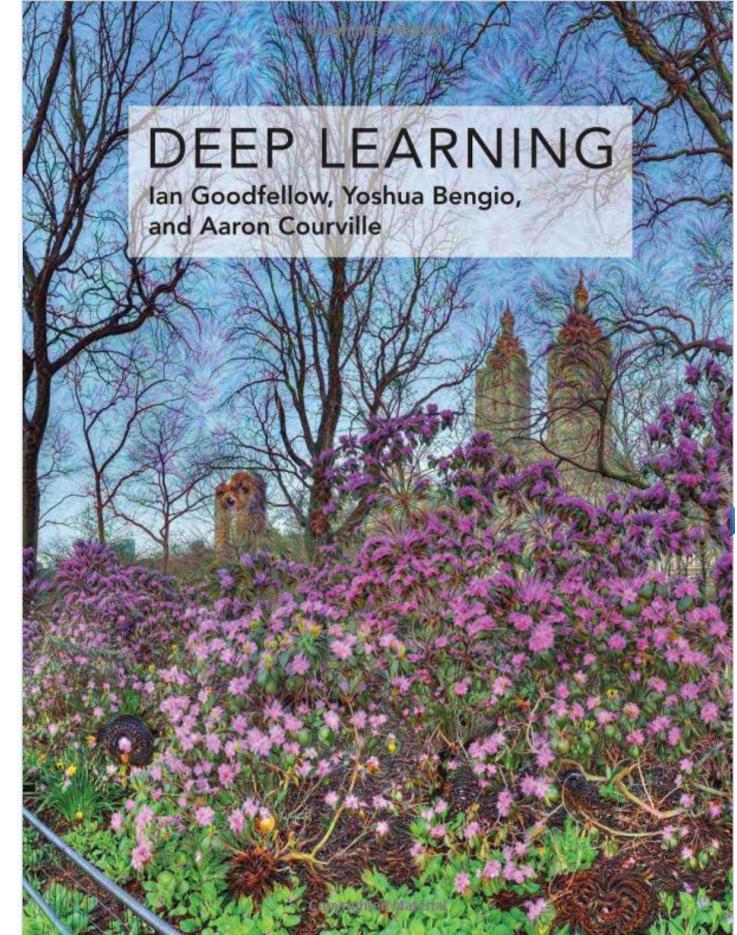


Lecture 2: Basics of machine learning

Reading:

Deep Learning, chapter 5 (up to and including section 5.5).

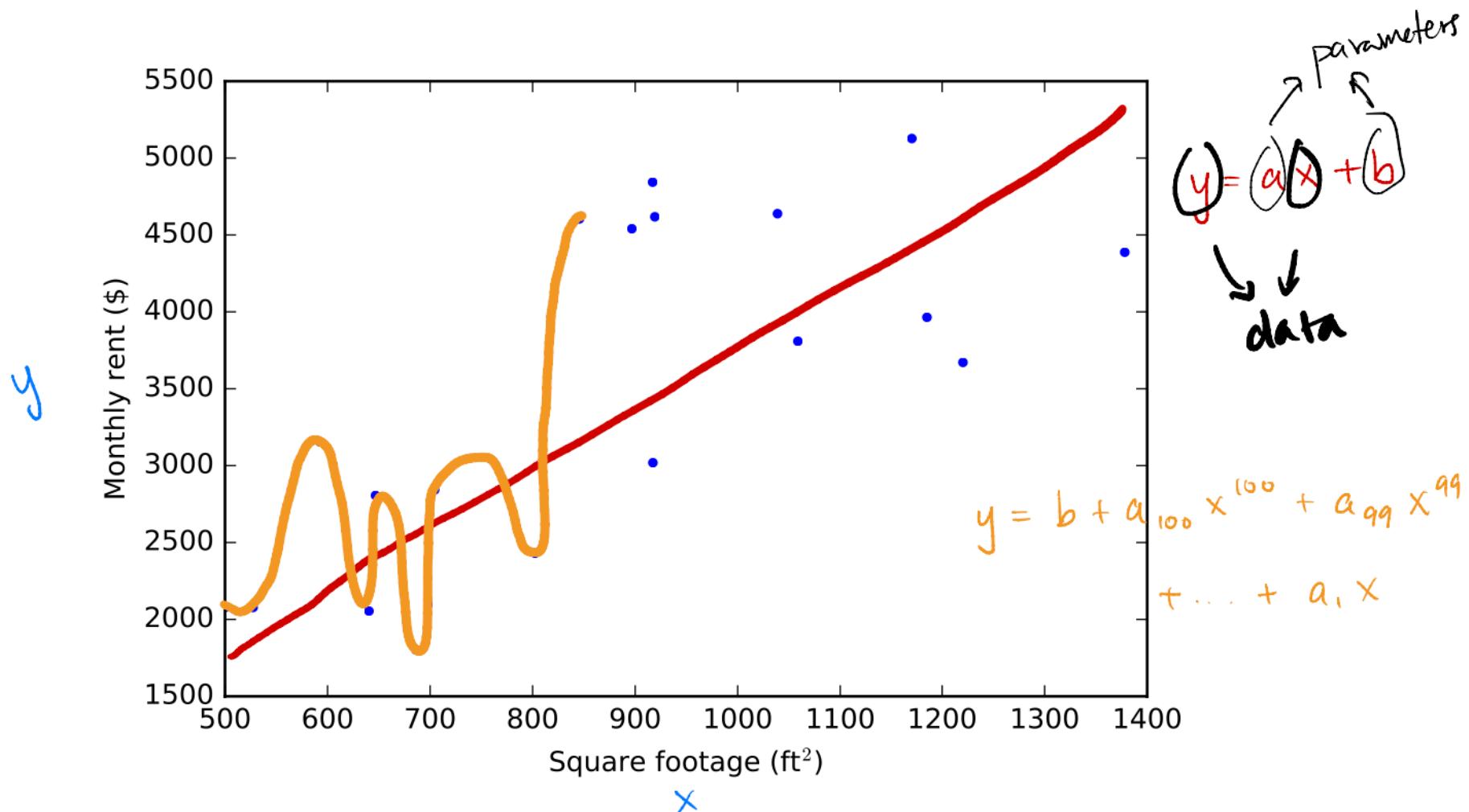
To refresh your linear algebra and probability, look at chapters 2 and 3 respectively.





An example of supervised learning

$$f(x) = ax + b$$

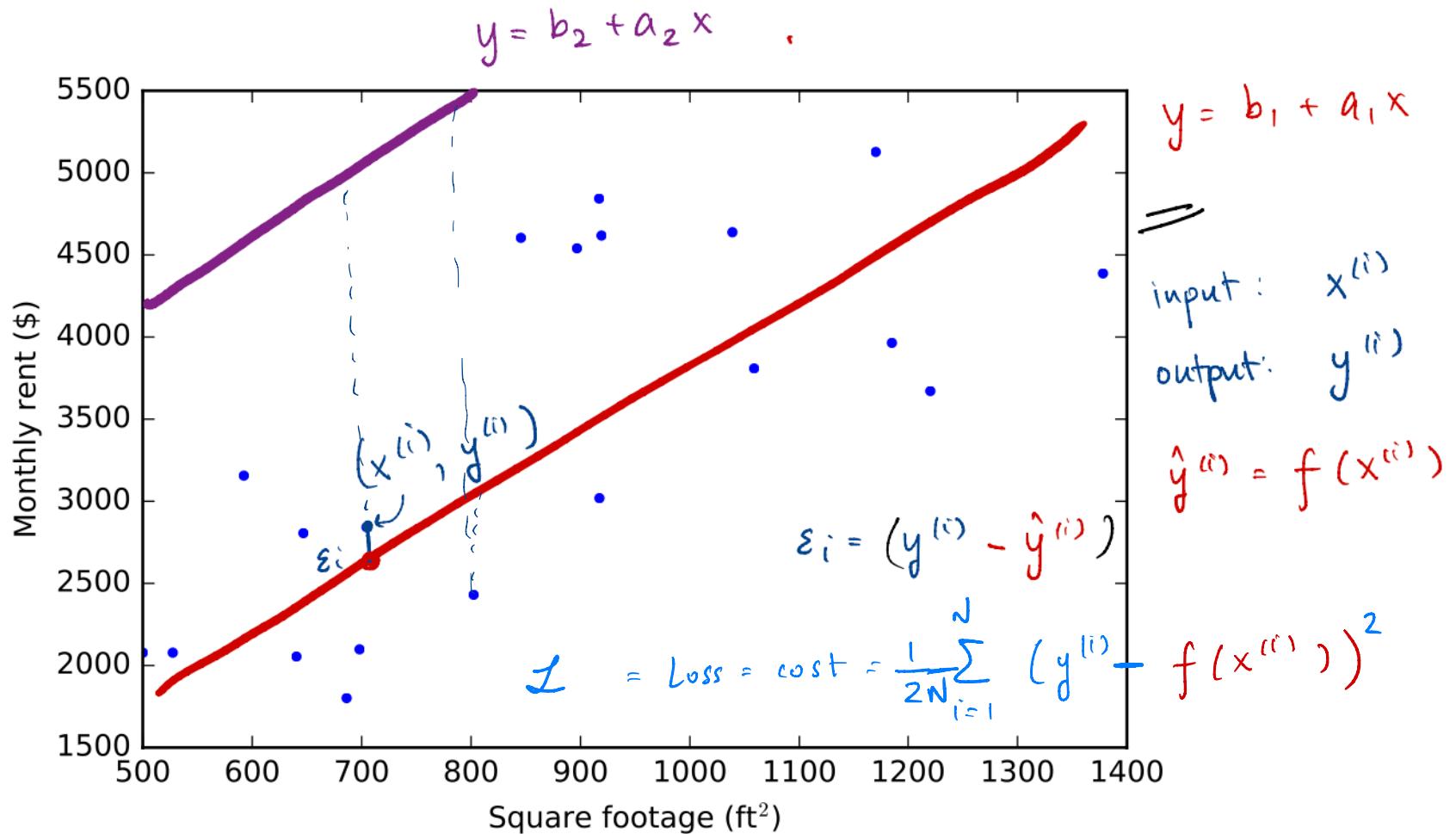


How should we model this data?

- ▶ Inputs, \mathbf{x} ? Outputs, \mathbf{y} ?
- ▶ What model should we use?
- ▶ How do we assess how good our model is?



An example of supervised learning



How should we model this data?

- ▶ Inputs, \mathbf{x} ? Outputs, \mathbf{y} ?
- ▶ What model should we use?
- ▶ How do we assess how good our model is?



An example of supervised learning

data given to us

A simple linear example:

- Model:

$$\hat{y} = ax + b \\ = \theta^T \hat{x}$$

parameters: I get to choose the values of a, b to make the model as good as poss.

- Cost function:

$$\theta = \begin{bmatrix} a \\ b \end{bmatrix} \quad \hat{x} = \begin{bmatrix} x \\ 1 \end{bmatrix} \text{ to make } \mathcal{L}(a, b) \text{ as small as poss.}$$

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

How do we learn the parameters of this model?



An example of supervised learning

Construct it as an optimization problem.

Our goal is to choose the parameters to make the loss as small as possible.

θ

$L(\theta)$

$L(\theta)$

$\frac{\partial L}{\partial \theta}$

θ

Thus our strategy to find the best θ is to:

- Calculate $\frac{dL}{d\theta}$
- Solve for θ such that $\frac{\partial L}{\partial \theta} = 0$

$$\frac{dL}{d\theta}$$

$$\frac{\partial L}{\partial \theta} = 0$$

However, θ is a vector, so how do we take derivatives with respect to it?



Aside: vector and matrix derivatives

In machine learning, we often take derivatives with respect to vectors and matrices.

These are typically called *gradients*, and in this class we'll use the following notation to denote derivatives with respect to vectors and matrices.

In this class, we will use both the differentiation operator ∂ and ∇ to denote derivatives. If y is a scalar, and \mathbf{x} is a vector, then the gradient of y with respect to \mathbf{x} is denoted as both:

$$\frac{\partial y}{\partial \mathbf{x}} \text{ and } \nabla_{\mathbf{x}} y$$

The gradient of y with respect to \mathbf{x} is itself a vector with the same dimensionality as \mathbf{x} .



Aside: vector and matrix derivatives

$$y \text{ scalar} \quad x \text{ vector} \in \mathbb{R}^n \quad \frac{\partial y}{\partial x} \text{ or } \nabla_x y \in \mathbb{R}^n$$

The gradient

The gradient generalizes the scalar derivative to multiple dimensions. Suppose $f : \mathbb{R}^n \rightarrow \mathbb{R}$ transforms a vector $x \in \mathbb{R}^n$ to a scalar. If $y = f(x)$, then the gradient is:

$$\nabla_x y = \begin{bmatrix} \frac{\partial y}{\partial x_1} \\ \frac{\partial y}{\partial x_2} \\ \vdots \\ \frac{\partial y}{\partial x_n} \end{bmatrix}$$

$$x = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

$$\frac{\partial y}{\partial x} = \begin{bmatrix} 1 \\ 0.5 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

$$\frac{\partial y}{\partial x_1} = 1$$

$$\frac{\partial y}{\partial x_2} = 0.5$$

$$\Delta y \text{ due to } \Delta x_1 \text{ is } \approx \frac{\partial y}{\partial x_1} \cdot \Delta x_1$$



Aside: vector and matrix derivatives

In other words, the gradient is:

- A vector that is the same size as \mathbf{x} , i.e., if $\mathbf{x} \in \mathbb{R}^n$ then $\nabla_{\mathbf{x}}y \in \mathbb{R}^n$.
- Each dimension of $\nabla_{\mathbf{x}}y$ tells us how small changes in \mathbf{x} in that dimension affect y . i.e., changing the i th dimension of \mathbf{x} by a small amount, Δx_i , will change y by

$$\frac{\partial y}{\partial x_i} \Delta x_i$$

We may also denote this as:

$$\Delta \vec{x} = \begin{bmatrix} 0.05 \\ 0.01 \\ 0.02 \\ \vdots \end{bmatrix} \quad (\nabla_{\mathbf{x}}y)_i \Delta x_i \rightarrow \text{how much does } y \text{ change?}$$

$$\Delta y \approx \sum_{i=1}^n \frac{\partial y}{\partial x_i} \Delta x_i$$

$$\approx (\nabla_{\mathbf{x}}y)^T \Delta \mathbf{x}$$



Aside: vector and matrix derivatives

Example: derivative with respect to a vector

Let $f(\mathbf{x}) = \theta^T \mathbf{x}$. What is $\nabla_{\mathbf{x}} f(\mathbf{x})$?

$$\theta, \vec{x} \in \mathbb{R}^n$$

$$y = f(x)$$

$$y \in \mathbb{R}$$

$$\nabla_{\mathbf{x}} y \in \mathbb{R}^n$$

$$y = \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$$

$$\nabla_{\mathbf{x}} y = \begin{bmatrix} \frac{\partial y}{\partial x_1} \\ \frac{\partial y}{\partial x_2} \\ \vdots \\ \frac{\partial y}{\partial x_n} \end{bmatrix} = \begin{bmatrix} \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{bmatrix}$$

$$\boxed{\nabla_{\mathbf{x}} f(\mathbf{x}) = \theta}$$



Aside: vector and matrix derivatives

Example: derivative with respect to a vector

Let $f(\mathbf{x}) = \mathbf{x}^T \mathbf{A} \mathbf{x}$. What is $\nabla_{\mathbf{x}} f(\mathbf{x})$?

$$\mathbf{x} \in \mathbb{R}^n \quad \mathbf{A} \in \mathbb{R}^{n \times n}$$

$(1 \times n) (n \times n) (n \times 1)$

$$\mathbf{A} = \begin{bmatrix} a_{11} & \dots & a_{1n} \\ \vdots & \ddots & \vdots \\ a_{n1} & \dots & a_{nn} \end{bmatrix}$$

$$\nabla_{\mathbf{x}} f(\mathbf{x}) \in \mathbb{R}^n$$

$i=1, j=1$

$$\frac{\partial a_{11} \cdot x_1^2}{\partial x_1} = 2a_{11} \cdot x_1$$

$i=1, j \neq 1$

$$f(\mathbf{x}) = \mathbf{x}^T \mathbf{A} \mathbf{x}$$
$$= \sum_{i=1}^n \sum_{j=1}^n a_{ij} \cdot x_i \cdot x_j$$

$$\frac{\partial \left(\sum_{j=2}^n a_{ij} \cdot x_i \cdot x_j \right)}{\partial x_1}$$

$$\frac{\partial f(\mathbf{x})}{\partial x_1} = 2a_{11}x_1 + \sum_{j=2}^n a_{1j} \cdot x_j + \sum_{i=2}^n a_{ii}x_i$$

$$= \sum_{j=1}^n a_{1j} \cdot x_j + \sum_{i=1}^n a_{ii}x_i$$

$i \neq 1, j=1$

$$\frac{\partial \sum_{i=2}^n a_{ii} \cdot x_i \cdot x_1}{\partial x_1}$$



Aside: vector and matrix derivatives

$$\frac{\partial f(x)}{\partial x_1} = \sum_{j=1}^n a_{1j} \cdot x_j + \sum_{i=1}^n a_{i1} \cdot x_i$$
$$(Ax)_1 + (A^T x)_1$$

$$\left[\begin{array}{c} \frac{\partial f(x)}{\partial x_1} \\ \frac{\partial f(x)}{\partial x_2} \\ \vdots \\ \frac{\partial f(x)}{\partial x_n} \end{array} \right] = \left[\begin{array}{c} (Ax)_1 + (A^T x)_1 \\ (Ax)_2 + (A^T x)_2 \\ \vdots \\ (Ax)_n + (A^T x)_n \end{array} \right]$$

$\frac{\partial f(x)}{\partial x_2}, \frac{\partial f(x)}{\partial x_3}$ all follow the same pattern

$$\nabla_x f(x) = \nabla_x (x^T A x) = (Ax + A^T x)$$

$$\begin{aligned}
 &= (A + A^T)x \\
 &= 2Ax
 \end{aligned}$$

if A is symmetric

when $n=1$

$$f(x) = x \cdot a \cdot x = a \cdot x^2$$

$$\frac{\partial f(x)}{\partial x} = 2ax$$



Aside: vector and matrix derivatives

First, we note that $\mathbf{x}^T \mathbf{A} \mathbf{x} = \sum_i \sum_j a_{ij} x_i x_j$. Then, we have

$$\begin{aligned}\nabla_{\mathbf{x}} f(\mathbf{x}) &= \begin{bmatrix} 2a_{11}x_1 + \color{red}{a_{12}x_2 + \cdots + a_{1n}x_n} + a_{21}x_2 + \cdots + a_{n1}x_n \\ 2a_{22}x_2 + \color{red}{a_{21}x_1 + \cdots + a_{2n}x_n} + a_{12}x_1 + \cdots + a_{n2}x_n \\ \vdots \\ 2a_{nn}x_n + \color{red}{a_{n1}x_1 + \cdots + a_{n,n-1}x_{n-1}} + a_{1n}x_1 + \cdots + a_{n-1,n}x_{n-1} \end{bmatrix} \\ &= \color{red}{\mathbf{Ax}} + \color{blue}{\mathbf{A}^T \mathbf{x}}\end{aligned}$$



Matrix derivatives

$$y = z^T A x$$

$$z \in \mathbb{R}^m$$

$$x \in \mathbb{R}^n$$

$$A \in \mathbb{R}^{m \times n}$$

$$\frac{\partial y}{\partial A} = \nabla_A y = \begin{bmatrix} \frac{\partial y}{\partial a_{11}} & \frac{\partial y}{\partial a_{12}} & \dots & \frac{\partial y}{\partial a_{1n}} \\ \vdots & \ddots & \ddots & \vdots \\ \frac{\partial y}{\partial a_{m1}} & \dots & \ddots & \frac{\partial y}{\partial a_{mn}} \end{bmatrix}$$

if $A \in \mathbb{R}^{m \times n}$

$$\nabla_A y \in \mathbb{R}^{m \times n}$$

"Denominator Layout"

"Numerator Layout"

$$x \in \mathbb{R}^{n \times 1}, \quad \nabla_x y \in \mathbb{R}^{n \times 1}$$
$$A \in \mathbb{R}^{m \times n}, \quad \nabla_A y \in \mathbb{R}^{m \times n}$$

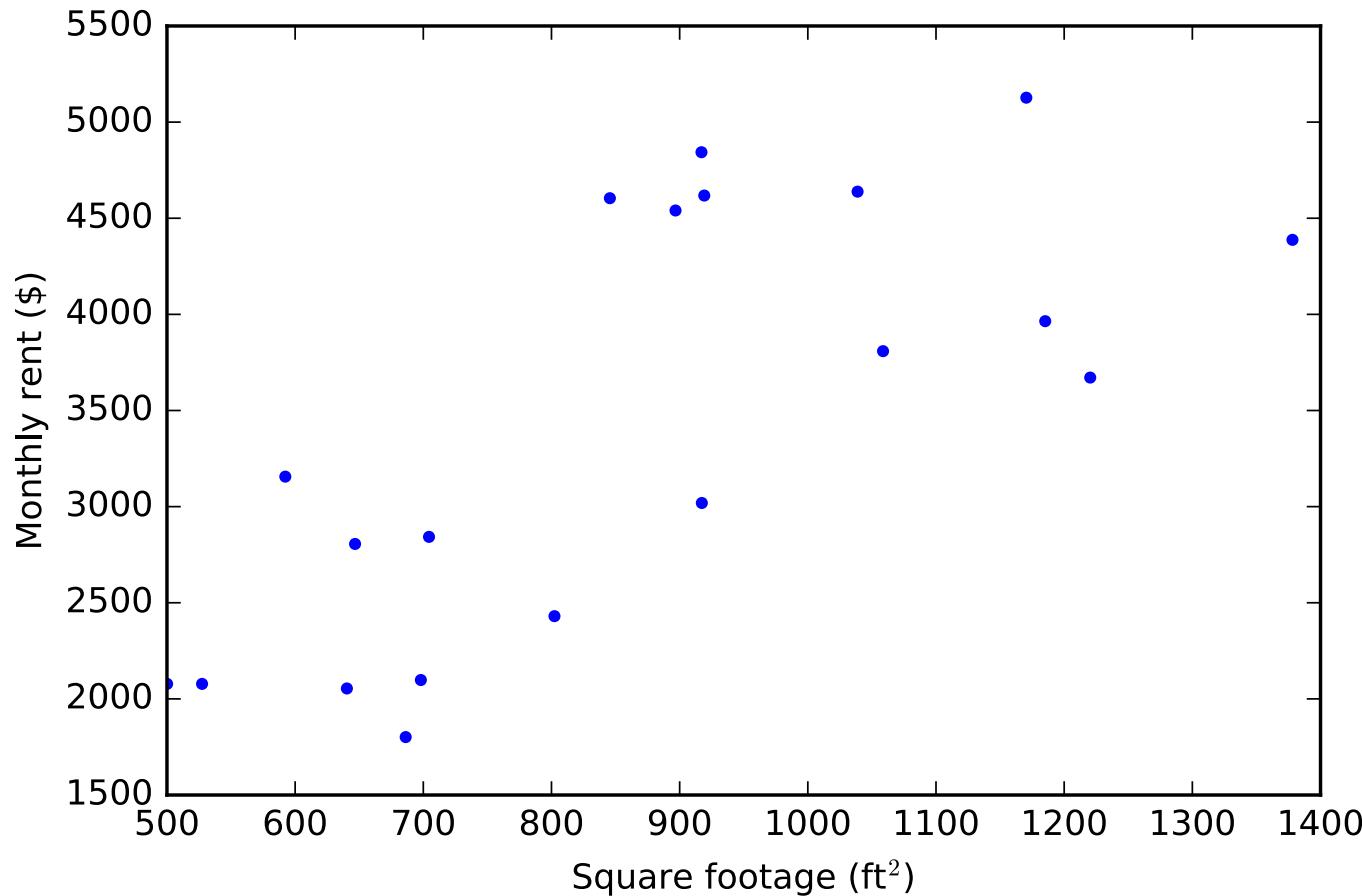
$$\nabla_x y \in \mathbb{R}^{1 \times n}$$
$$\nabla_A y \in \mathbb{R}^{n \times m}$$

In numerator layout

$$\frac{\partial y}{\partial A} = \begin{bmatrix} \frac{\partial y}{\partial a_{11}} & \frac{\partial y}{\partial a_{21}} & \dots \\ \frac{\partial y}{\partial a_{12}} \\ \vdots \\ \frac{\partial y}{\partial a_{1n}} \end{bmatrix}$$



Back to our supervised learning example



Thus our strategy to find the best θ is to:

- Calculate $\frac{d\mathcal{L}}{d\theta}$
- Solve for θ such that $\frac{\partial \mathcal{L}}{\partial \theta} = 0$

$$\frac{d\mathcal{L}}{d\theta}$$

$$\frac{\partial \mathcal{L}}{\partial \theta} = 0$$



Back to our supervised learning example

$$\theta = \begin{bmatrix} a \\ b \end{bmatrix} \quad \hat{x} = \begin{bmatrix} x \\ 1 \end{bmatrix}$$

Re-writing the cost function, we have:

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

$$\begin{aligned}
 & \cancel{\theta^T = \theta} \\
 & = \frac{1}{2} \sum_{i=1}^N (y^{(i)} - \theta^T \hat{x}^{(i)})^T (y^{(i)} - \theta^T \hat{x}^{(i)}) \\
 & = \frac{1}{2} \sum_{i=1}^N (y^{(i)} - \hat{x}^{(i)\top} \theta)^T (y^{(i)} - \hat{x}^{(i)\top} \theta) \\
 & \xleftarrow{\text{"vectorization"}}
 & = \frac{1}{2} \left[\begin{array}{c} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(N)} \end{array} \right] - \left[\begin{array}{c} \hat{x}^{(1)\top} \\ \hat{x}^{(2)\top} \\ \vdots \\ \hat{x}^{(N)\top} \end{array} \right] \theta^T \left[\begin{array}{c} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(N)} \end{array} \right] - \left[\begin{array}{c} \hat{x}^{(1)\top} \\ \hat{x}^{(2)\top} \\ \vdots \\ \hat{x}^{(N)\top} \end{array} \right] \theta \\
 & \quad \underbrace{Y \in \mathbb{R}^N}_{\mathcal{Y}} \quad \underbrace{X \in \mathbb{R}^{N \times 2}}_{\mathcal{X}} \quad \theta \in \mathbb{R}^2 \\
 & = \frac{1}{2} (\mathcal{Y} - \mathcal{X}\theta)^T (\mathcal{Y} - \mathcal{X}\theta) \approx \frac{1}{2} \|\mathcal{Y} - \mathcal{X}\theta\|_2^2 \\
 & = \frac{1}{2} (\mathcal{Y}^T \mathcal{Y} - \mathcal{Y}^T \mathcal{X}\theta - \theta^T \mathcal{X}^T \mathcal{Y} + \theta^T \mathcal{X}^T \mathcal{X}\theta)
 \end{aligned}$$



$$A = (BC)$$

$$A^T = C^T B^T$$

Back to our supervised learning example

$$\frac{\partial \theta^T A \theta}{\partial \theta} = (A + A^T)\theta$$

$$w = z^T \theta \quad \frac{\partial w}{\partial \theta} = z$$

Let's now take derivatives:

$$\begin{aligned}\mathcal{L}(\theta) &= \frac{1}{2} \left(Y^T Y - \underbrace{2Y^T \mathbf{X} \theta}_{\substack{(1 \times N) (N \times 2) \\ 1 \times 2}} + \theta^T \mathbf{X}^T \mathbf{X} \theta \right) \\ \frac{\partial \mathcal{L}}{\partial \theta} &= \frac{1}{2} \left[0 - 2 \mathbf{X}^T Y + (\mathbf{X}^T \mathbf{X} + \mathbf{X}^T \mathbf{X}) \theta \right]\end{aligned}$$

$$= -\mathbf{X}^T Y + \mathbf{X}^T \mathbf{X} \theta$$

$$[=] 0$$

$$\mathbf{X}^T Y = \mathbf{X}^T \mathbf{X} \theta$$

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



Back to our supervised learning example



Back to our supervised learning example

This solution is called least-squares, and it appears in a variety of linear applications.

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



Back to our supervised learning example

$y \quad (N,)$ away
 $x \quad (N,)$ away

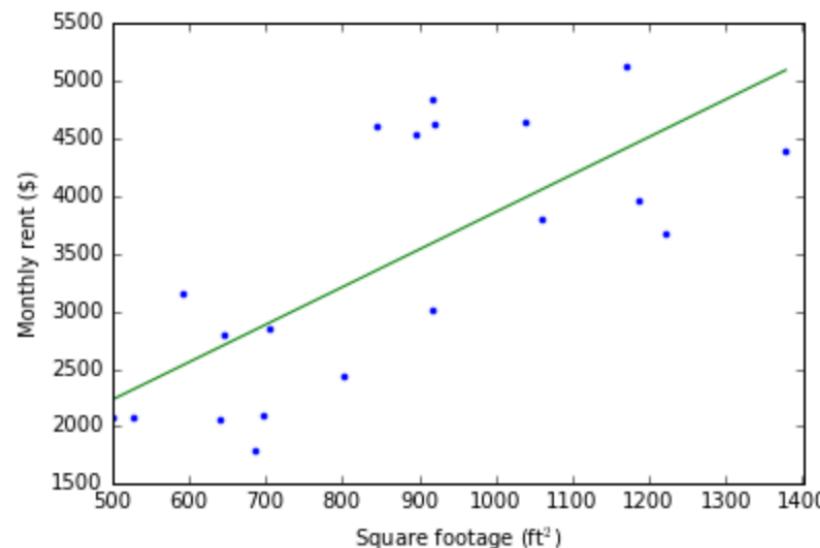
```
# Given paired data, with x being square footages and y being monthly rents
# Fit a linear model
xhat = np.vstack((x, np.ones_like(x)))
theta_lin = np.linalg.inv(xhat.dot(xhat.T)).dot(xhat.dot(y))

# Plot the data
f = plt.figure()
ax = f.gca()
ax.plot(x, y, '.')
ax.set_xlabel('Square footage (ft$^2$)')
ax.set_ylabel('Monthly rent ($)')
f.savefig('ml-basics_rent.pdf')

# Plot the regression line
xs = np.linspace(min(x), max(x), 50)
xs = np.vstack((xs, np.ones_like(xs)))
plt.plot(xs[0,:], theta_lin.dot(xs))

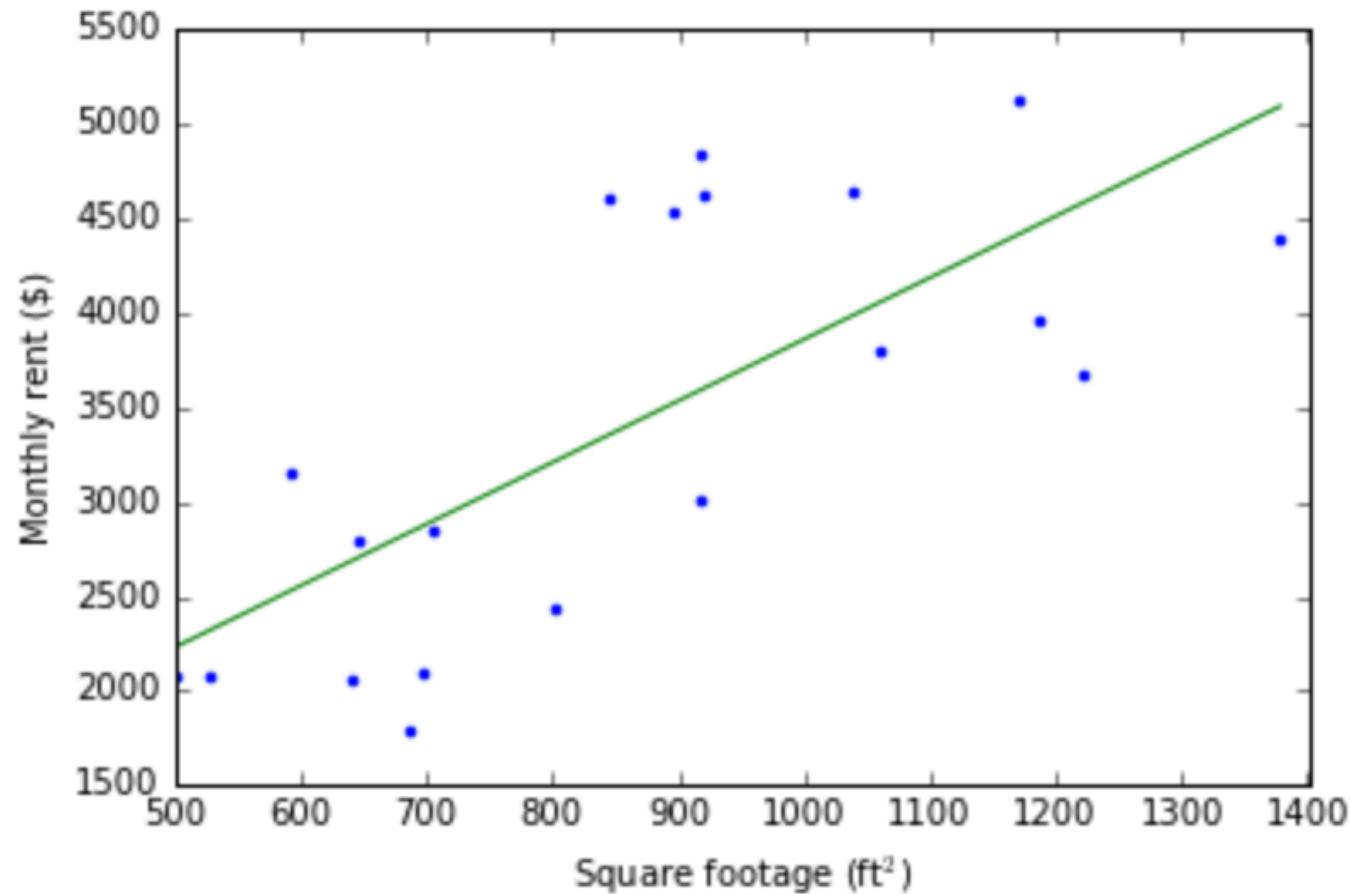
print theta_lin
```

```
[ 3.25900793  601.85544895]
```



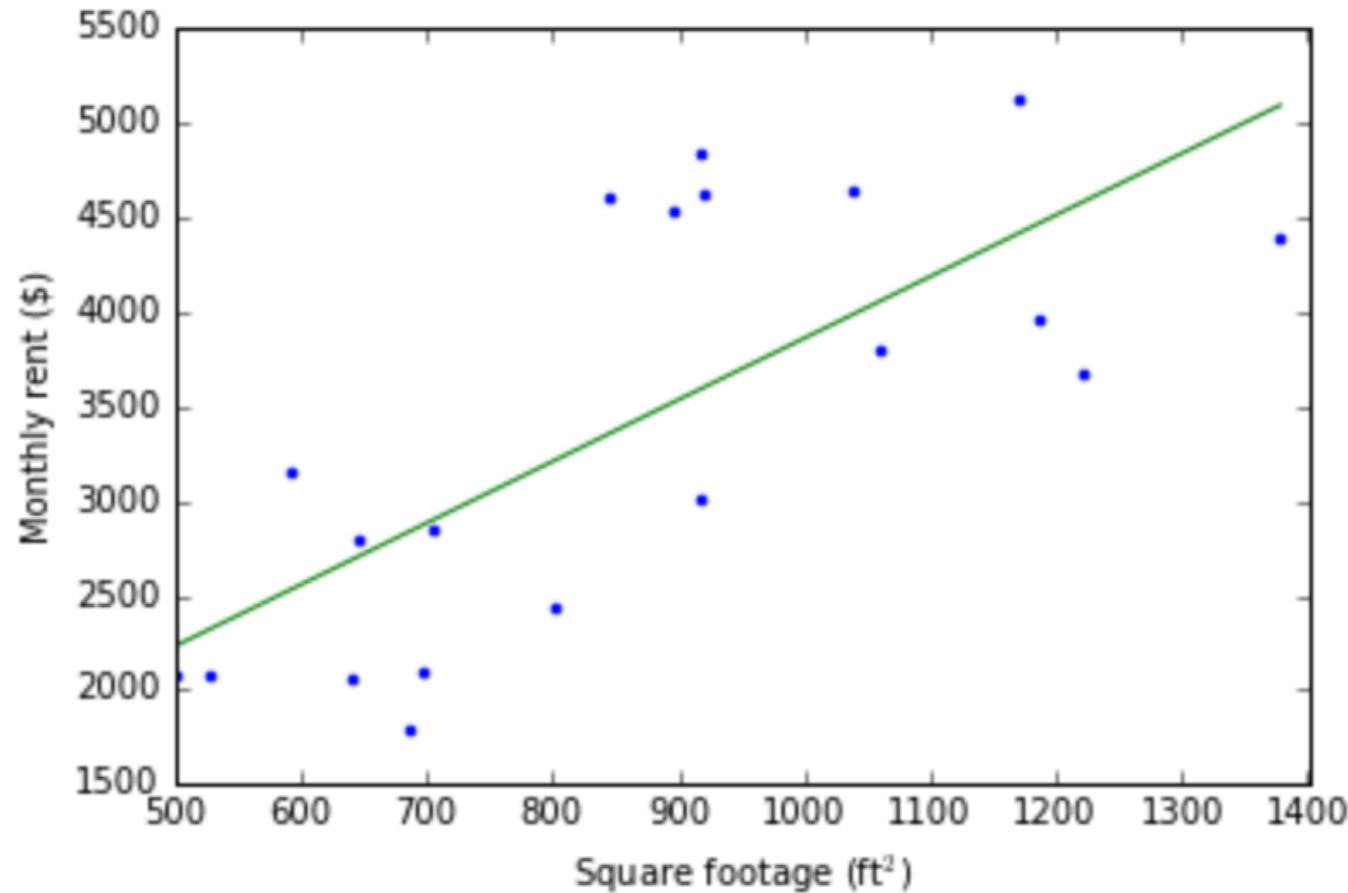


Can't we do better?



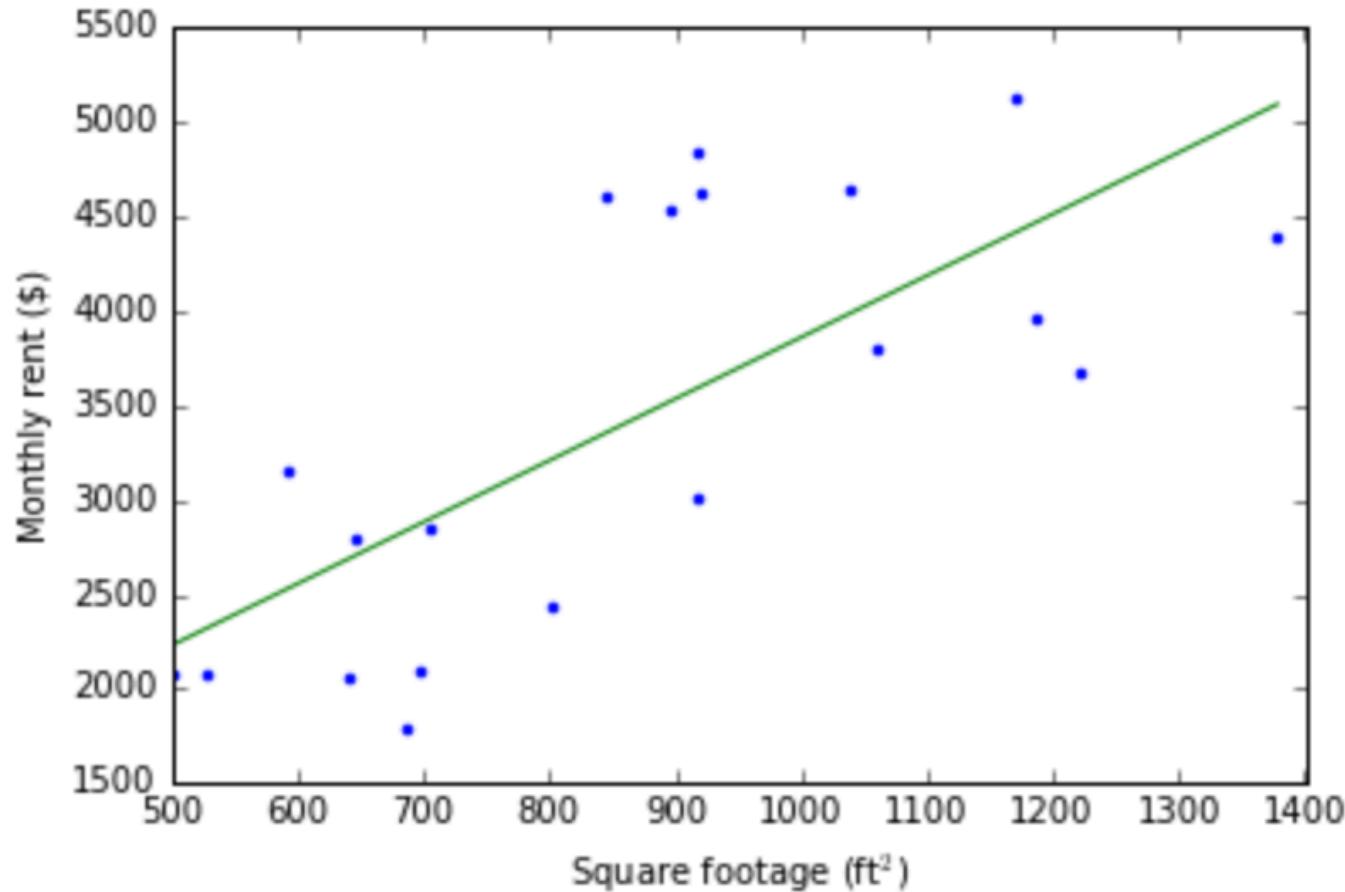


Can't we do better?





Question:



Question: How does our current least-squares formula allow for learning nonlinear polynomial fits, e.g.,

$$y = b + a_1x_1 + a_2x^2 + \cdots + a_nx^n$$

$$\hat{x} = \begin{bmatrix} x^3 \\ x^2 \\ x \\ 1 \end{bmatrix}$$



More generally...

With our problem setup, it's straightforward to generalize to higher degree polynomial models, e.g.,

$$y = b + a_1x + a_2x^2 + \cdots + a_nx^n$$

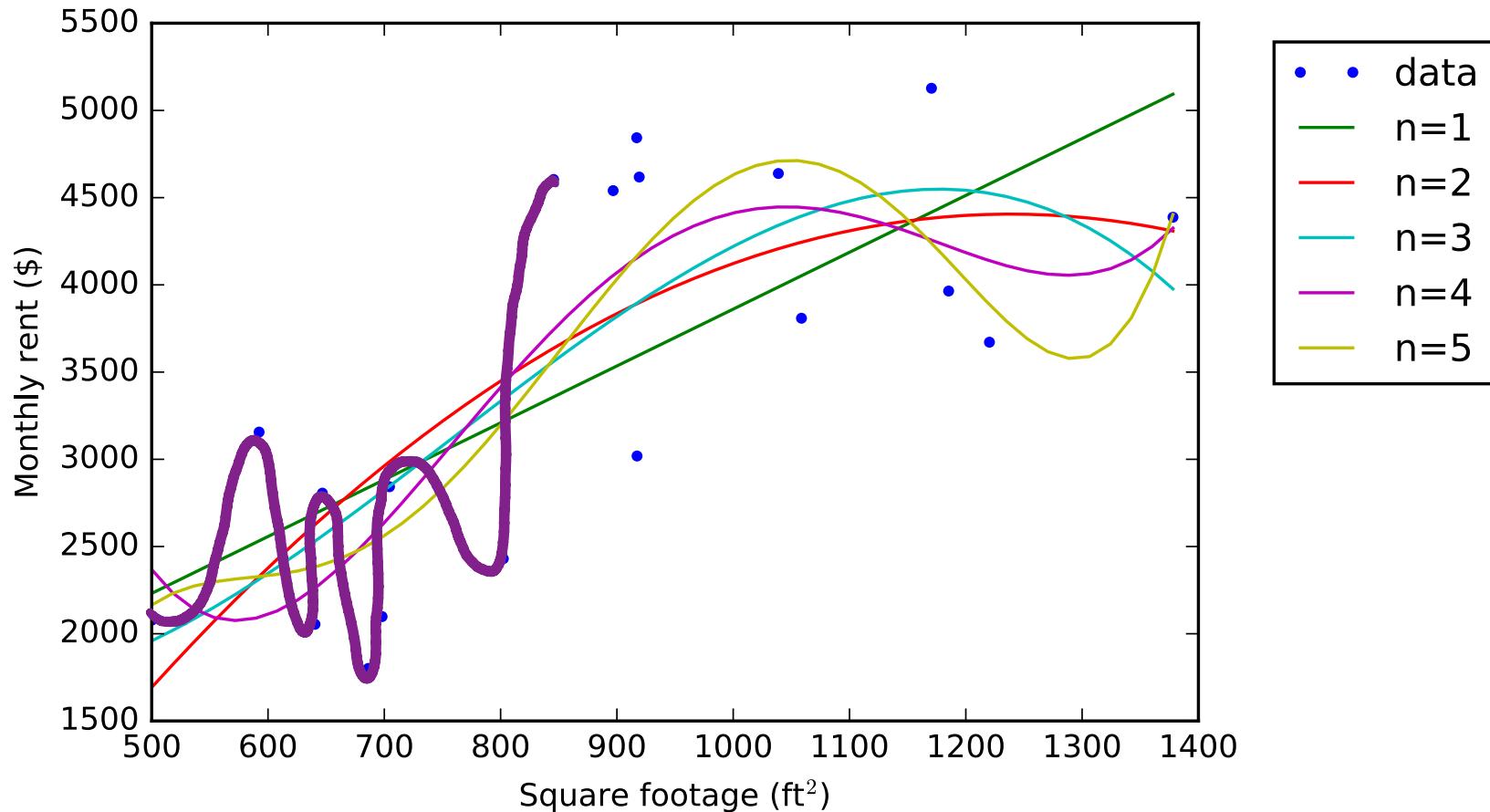
$$\hat{x} = \begin{bmatrix} 1 \\ x \\ x^2 \\ \vdots \\ x^n \end{bmatrix} \quad \theta = \begin{bmatrix} b \\ a_1 \\ a_2 \\ \vdots \\ a_n \end{bmatrix}$$

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



More generally...

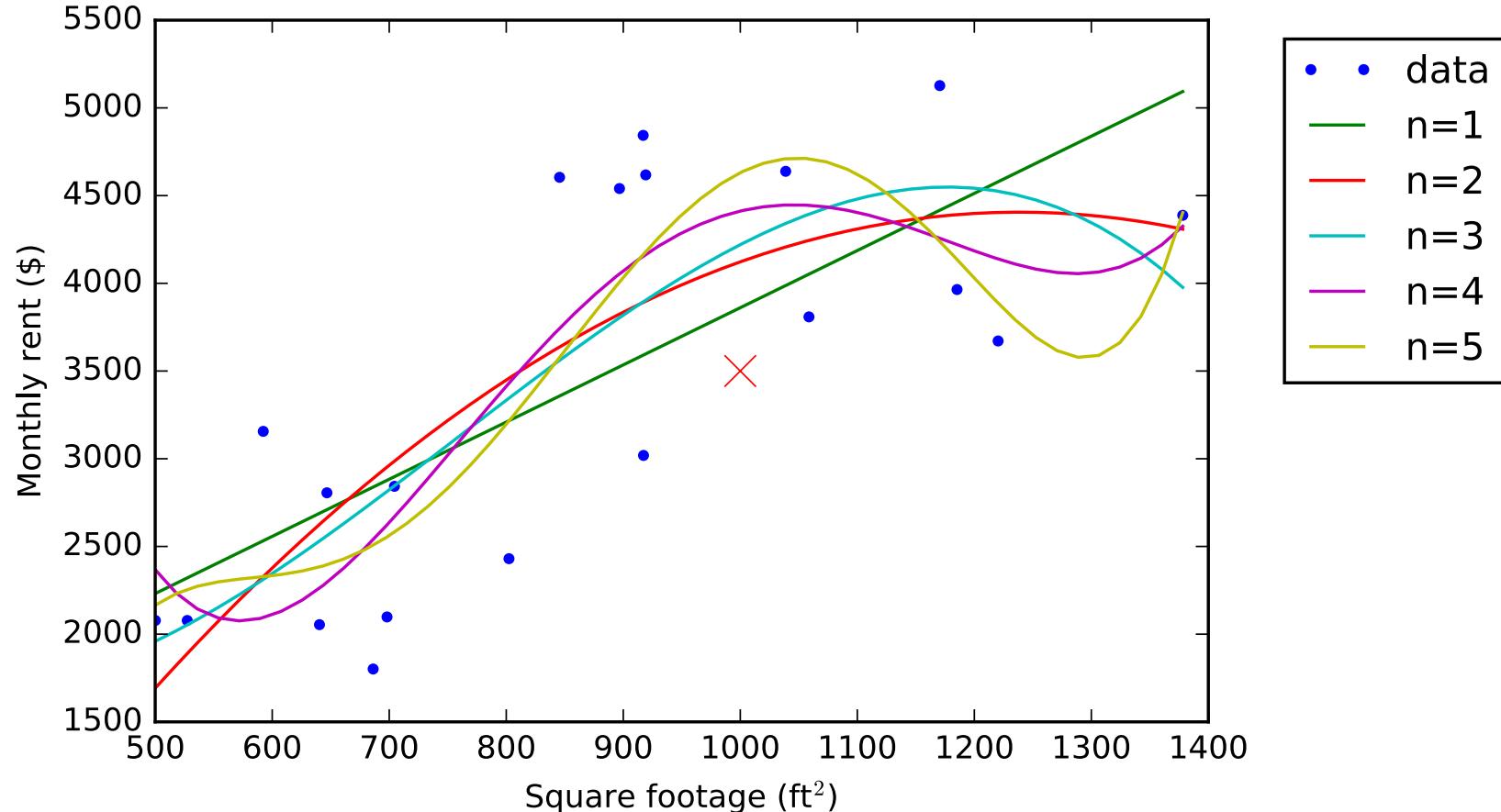
$$y = b + a_1 x + a_2 x^2 + a_3 x^3$$



Now, a higher degree polynomial will *always* fit the **provided** data as well as a lower order polynomial. Why?



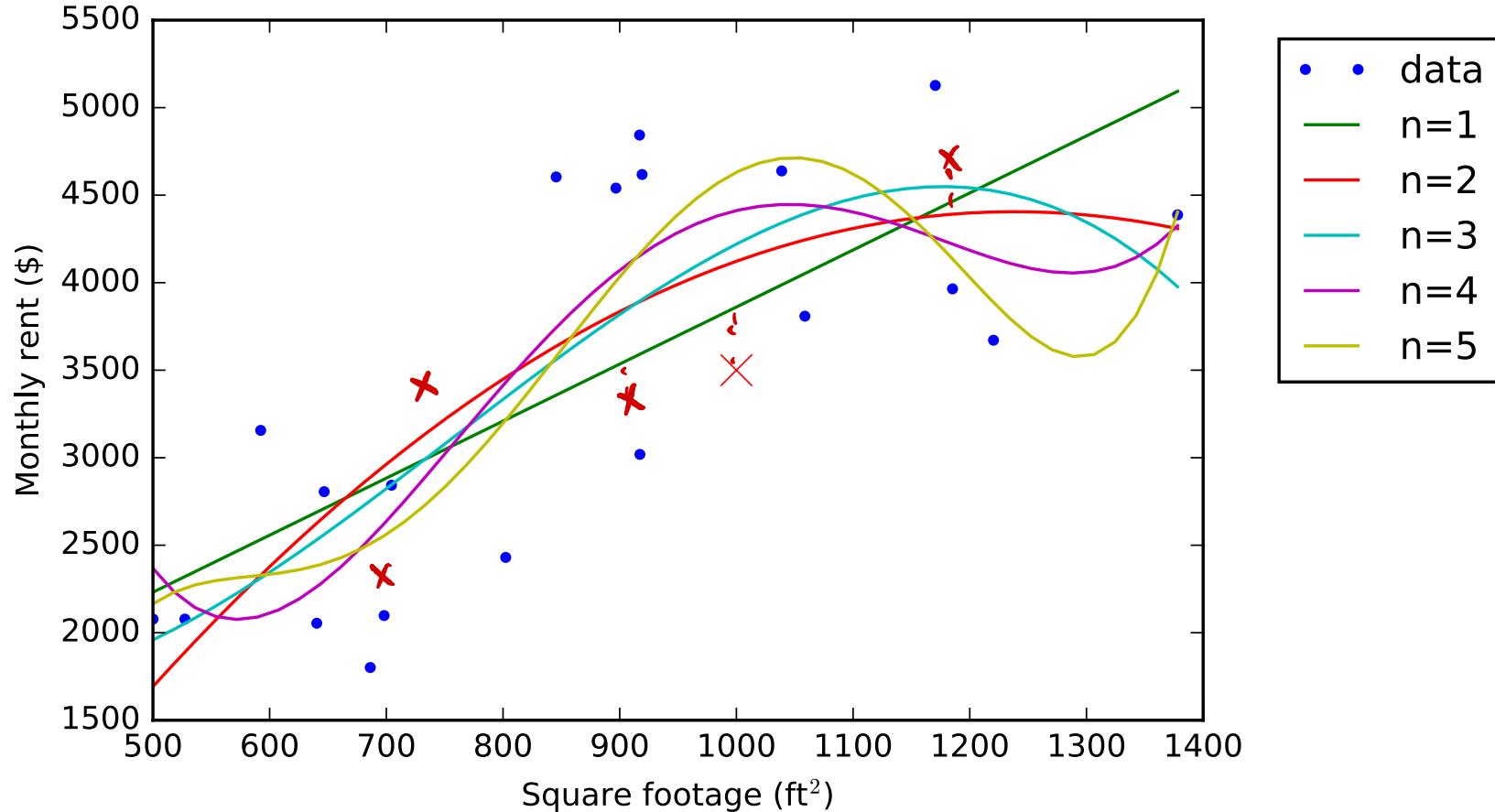
Now, let's say a new house pops up for rent...



In this scenario, the linear model best predicts the price of the new house ('X')



Model generalization



The fundamental problem here is that the more complex models may **not generalize as well** if the data come from a different model.



Training and testing data

Overfitting

This idea of generalization can be made more formal by introducing the concepts of a *training set* and *testing set*.

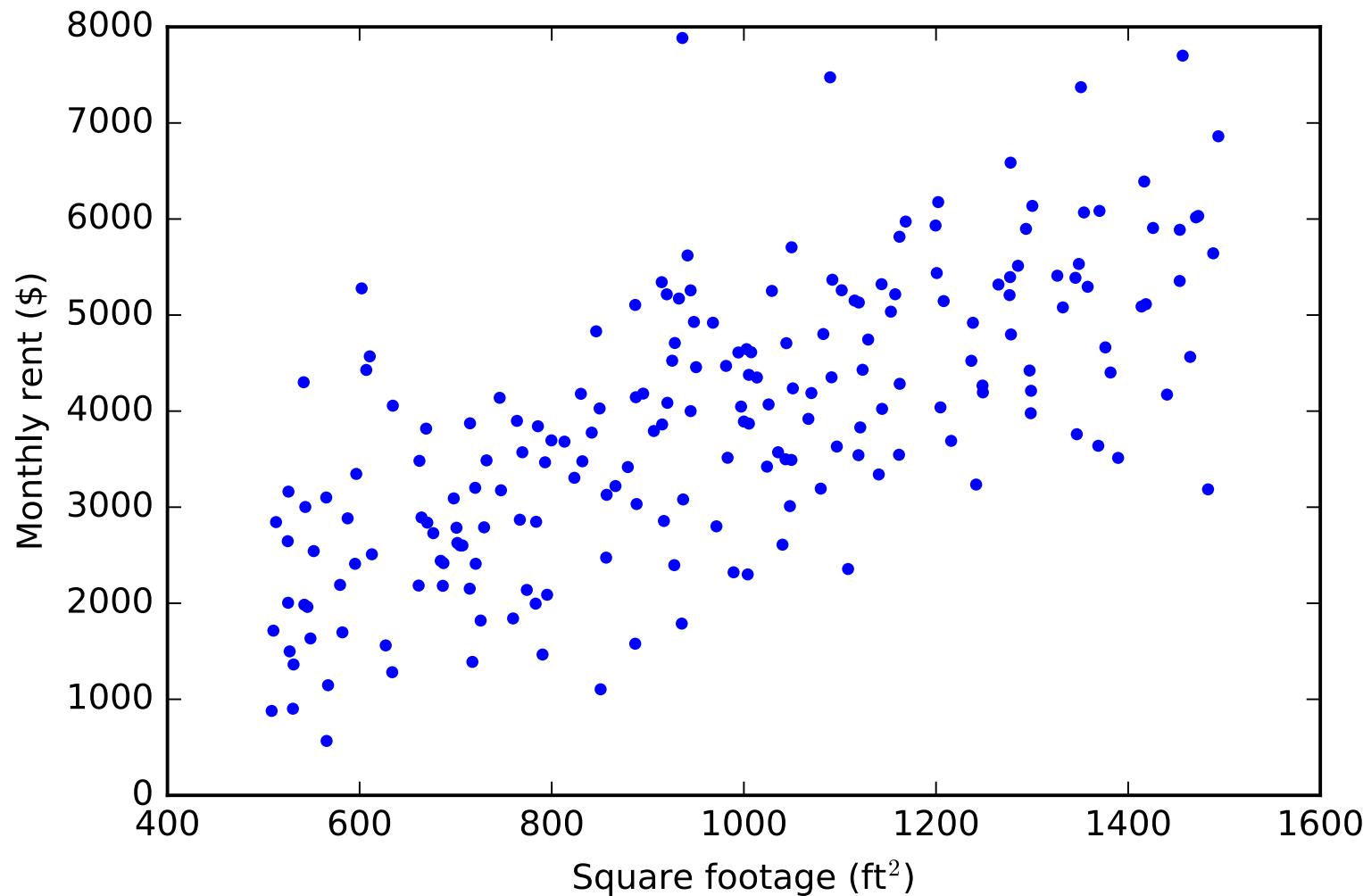
- **Training data** is data that is used to learn the parameters of your model.
- **Testing data** is data that is excluded in training and used to score your model.

There is also a notion of validation data, which we will get to later.

A model which has very low training error but high testing error is called *overfit*.

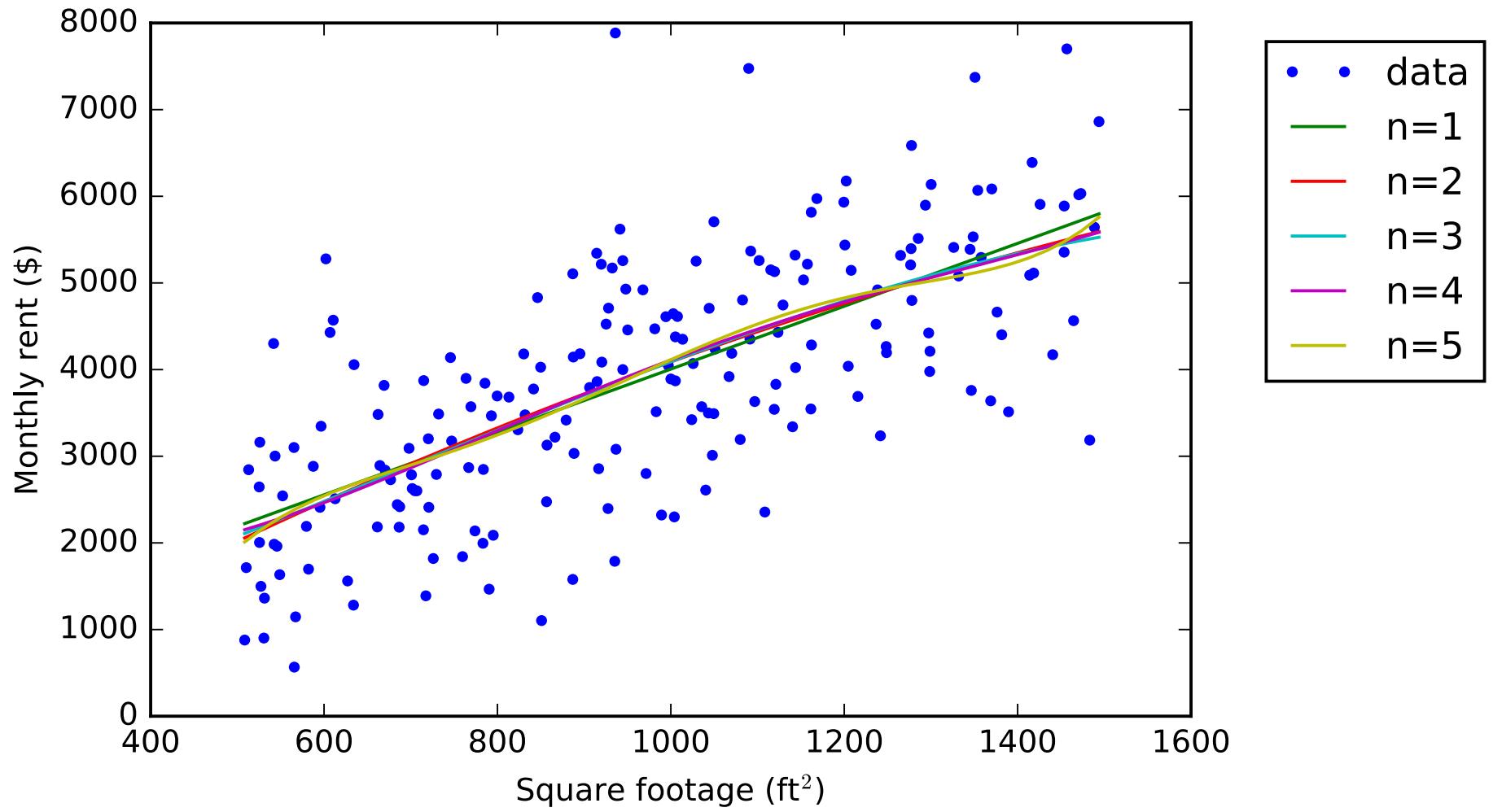


More data helps to avoid overfitting





More data helps to avoid overfitting





More data helps to avoid overfitting

This suggests that when *a lot* of data is available, it may be appropriate to use more complex models. (Another technique we will discuss later, regularization, also helps with overfitting.) This is a shadow of things to come (i.e., neural networks which have large complexity).



Underfitting

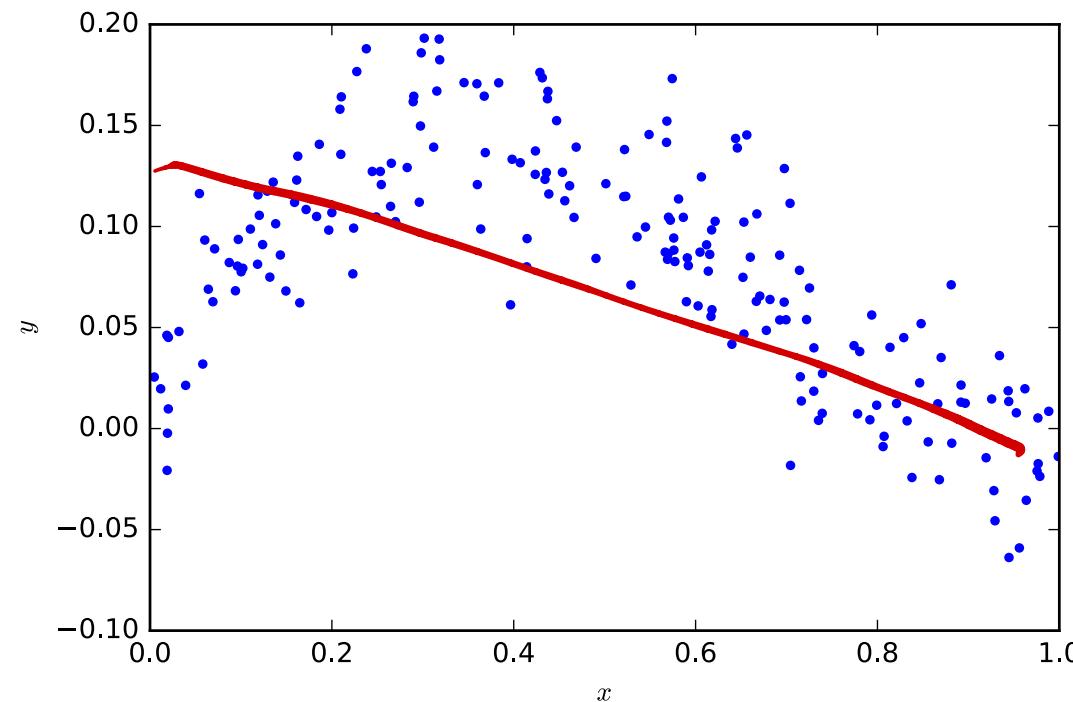
At the same time, we can not make the model *overly simple* as then we will underfit the data. This corresponds to a model that has both high training and high testing error, and the fit generally will not improve with more training data because the model is not expressive enough.



Underfitting

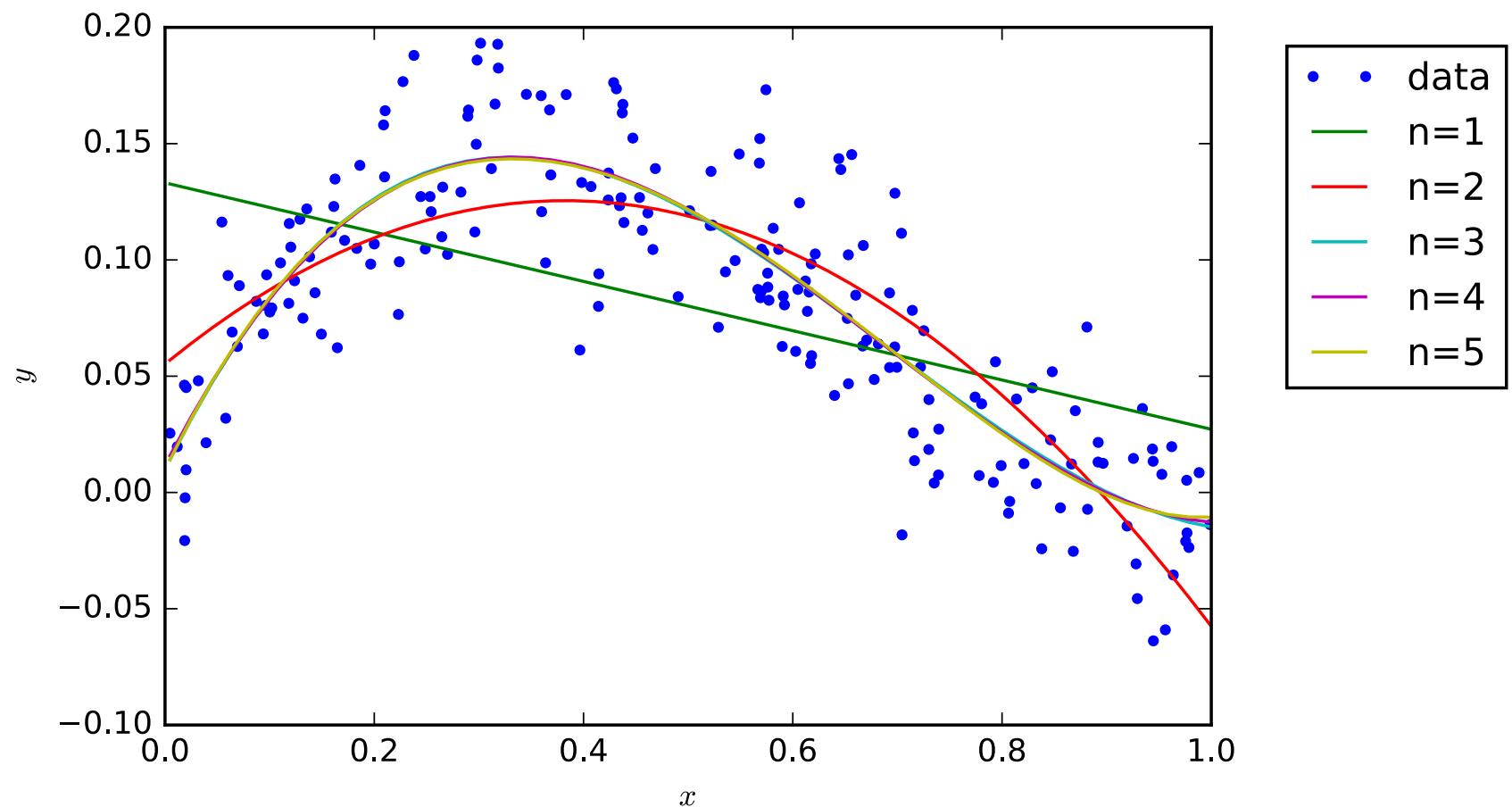
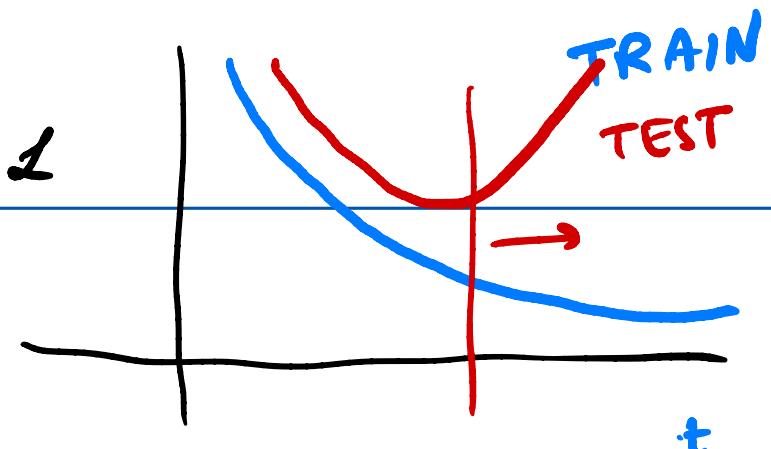
$$y = ax + b$$

```
np.random.seed(0) # Sets the random seed.  
num_train = 200 # Number of training data points  
  
# Generate the training data  
x = np.random.uniform(low=0, high=1, size=(num_train,))  
y = x - 2*x**2 + x**3 + np.random.normal(loc=0, scale=0.03, size=(num_train,))  
f = plt.figure()  
ax = f.gca()  
ax.plot(x, y, '.')  
ax.set_xlabel('$x$')  
ax.set_ylabel('$y$')  
f.savefig('ml-basics_underfitting-data.pdf')
```





Underfitting





Evaluating generalization error

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

4-fold
CV

This assumes you already have a *SEPARATE* testing set.
The TRAINING set is split into training and validation data.

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. **$N=800$**
- Then, split the data into k equal sets, each of N/k examples. Each of these sets is called a “fold.” *4 folds, each with 200 examples*
- $k - 1$ of the folds are datasets that are used to train the model parameters. *validation*
- 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

Original data

k-fold cross-validation with no hyperparameters

Very rare, and in this case you can think of this testing error as computing an average testing error of your model.



k-fold cross-validation with hyperparameters

The prevailing paradigm.





Evaluating generalization error



$$N = 1000$$

