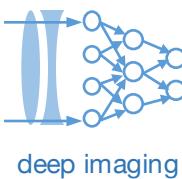


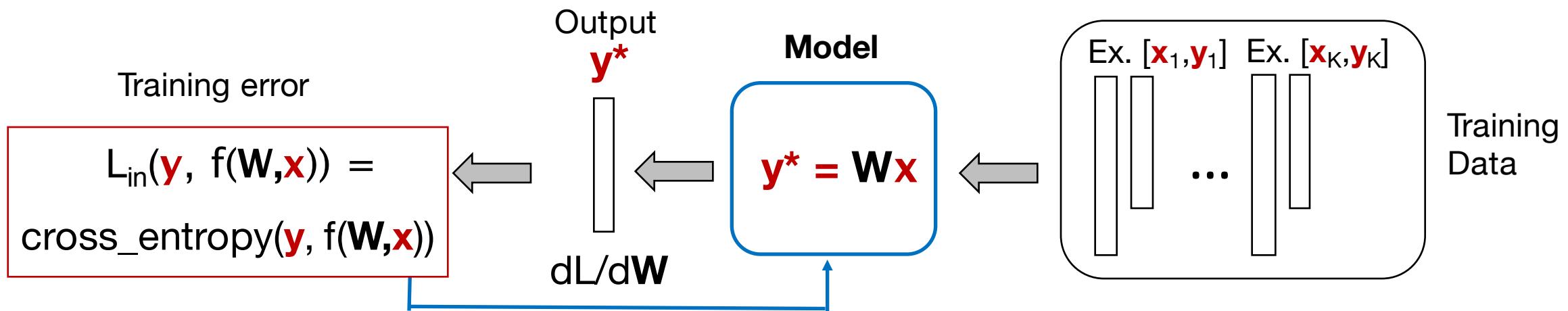
Lecture 8: Theoretical basics of machine learning

Machine Learning and Imaging

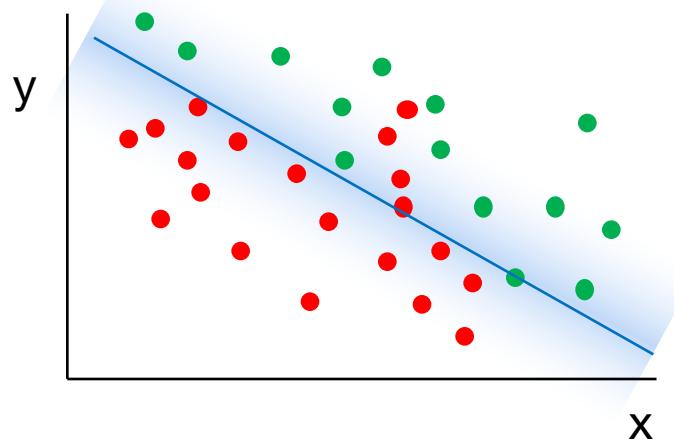
BME 590L
Roarke Horstmeyer



Review: the linear classification model – what's not to like?

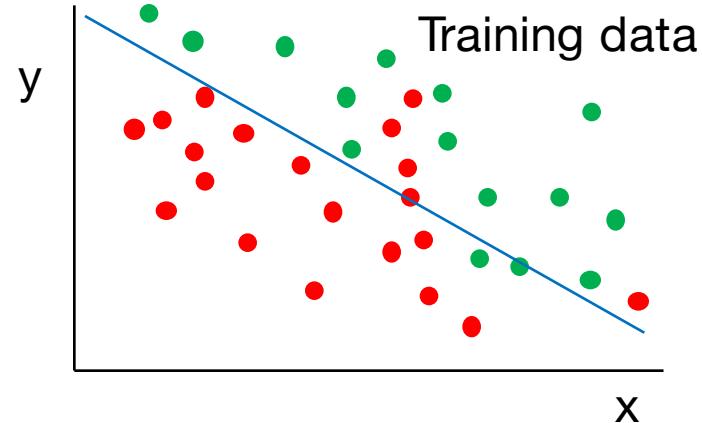
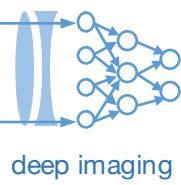


Probabilistic mapping to y



Limitations of linear classification model:

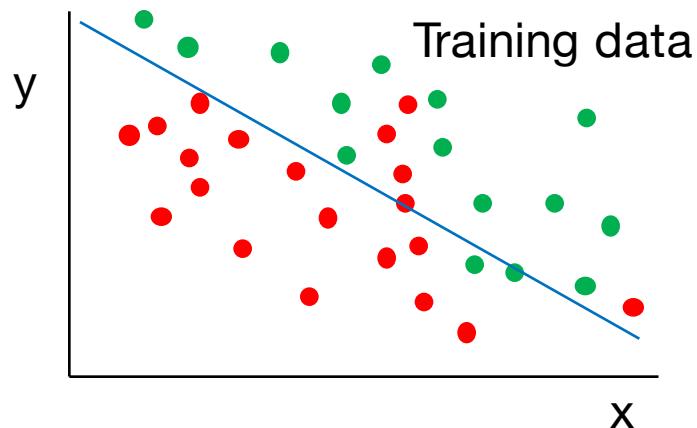
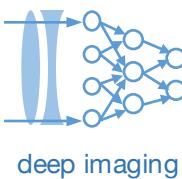
1. Can only separate data with lines (hyper-planes)...
2. We only allowed for binary labels ($y = +/- 1$)
3. Error function L_{in} inherently makes assumptions about statistical distribution of data



$$f = W_1 x$$

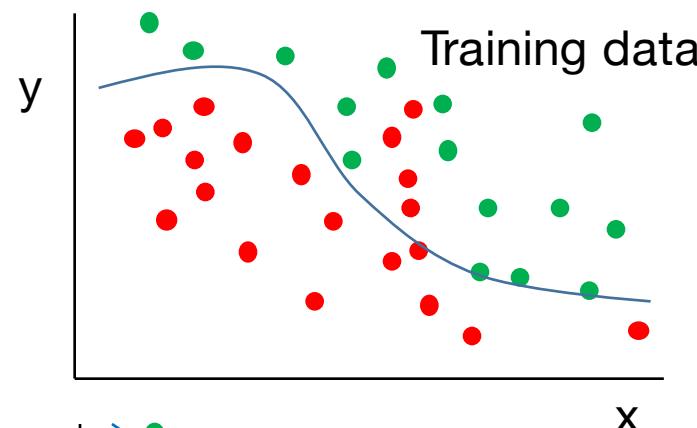
Learned f : not flexible

$$f = W_1 x$$



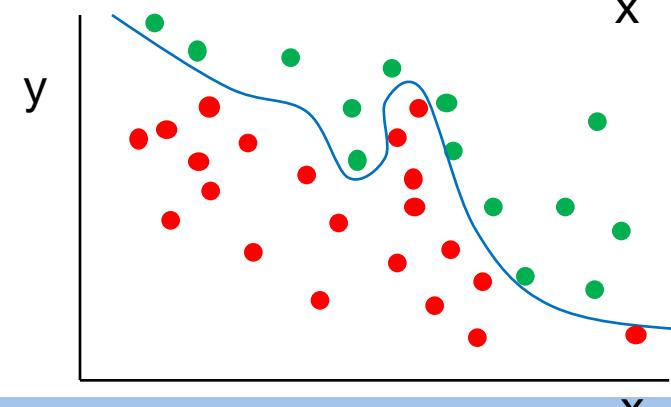
$$f = W_1 x$$

Learned f : not flexible



$$f = W_2 \max(W_1 x, 0)$$

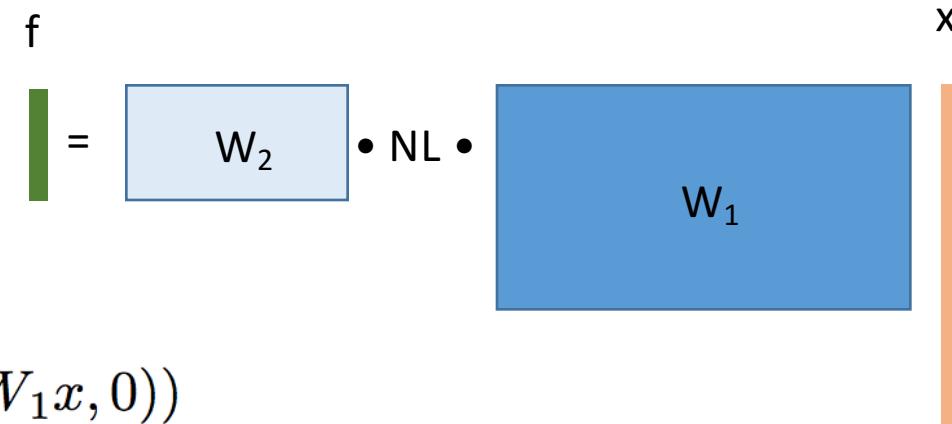
Learned f : a bit flexible



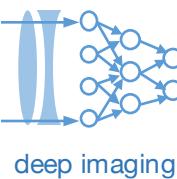
$$f = W_3 \max(0, W_2 \max(W_1 x, 0))$$

Learned f : more flexible

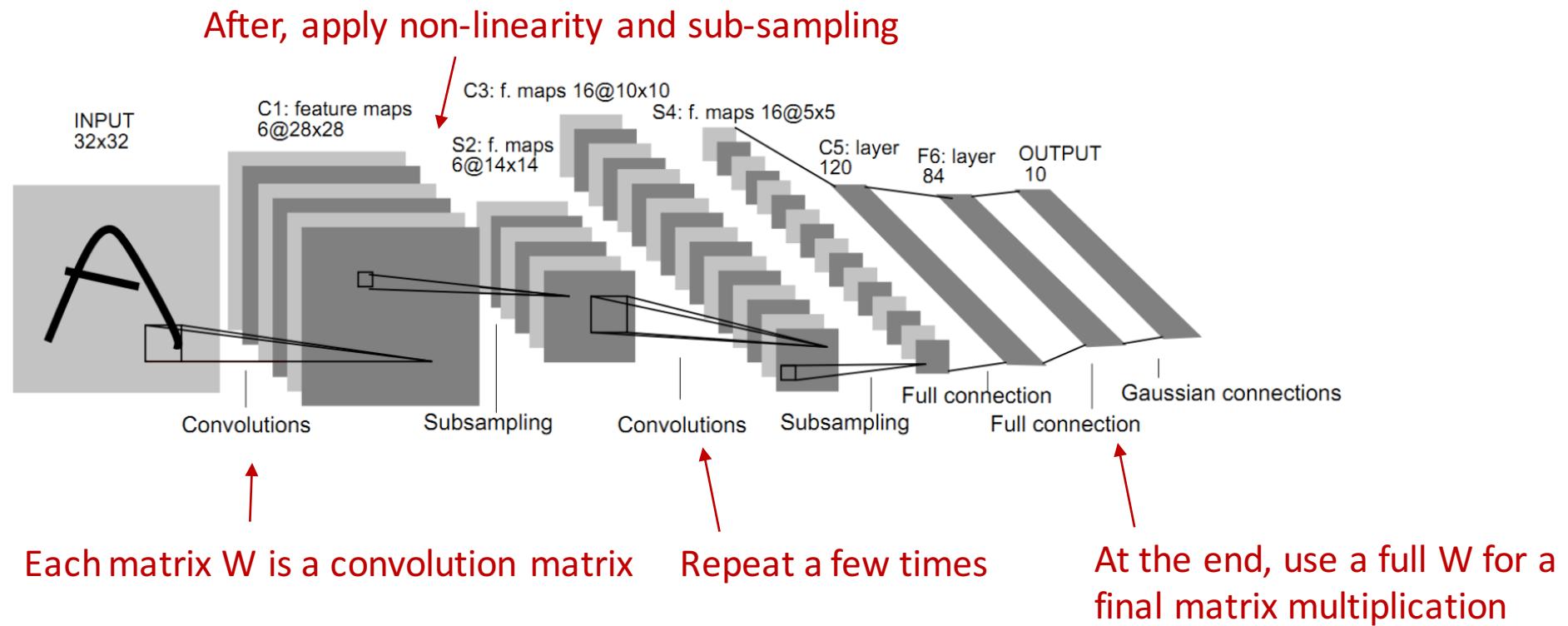
Does it generalize???



We can keep adding
these “layers”...

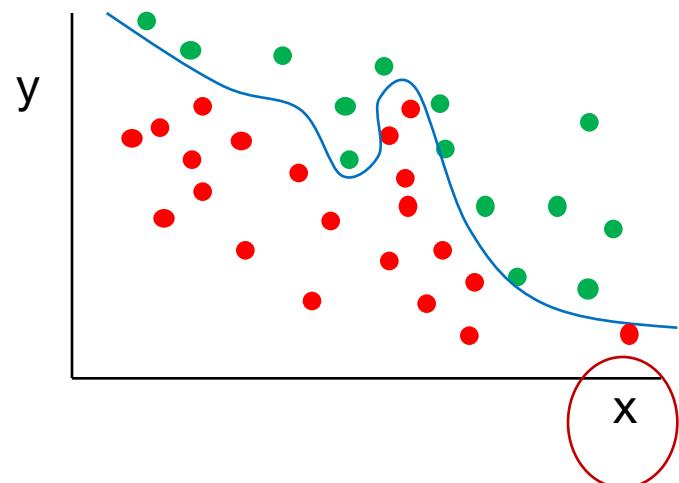


Getting us to Convolutional Neural Networks



Aside #1 before convolutional neural network details

Q: Can we try to avoid making these learning models too complicated?

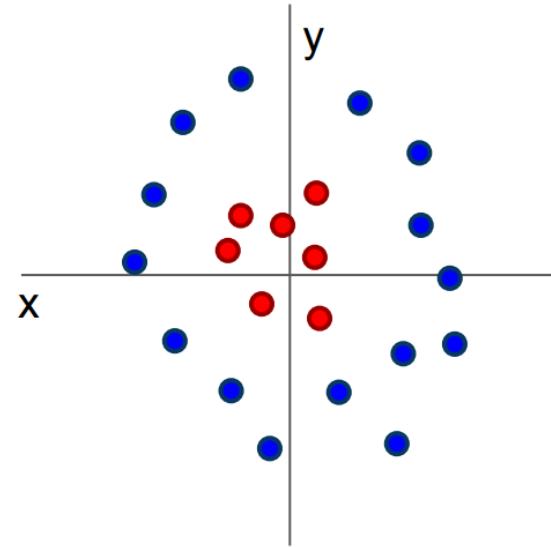


Learned f : more flexible

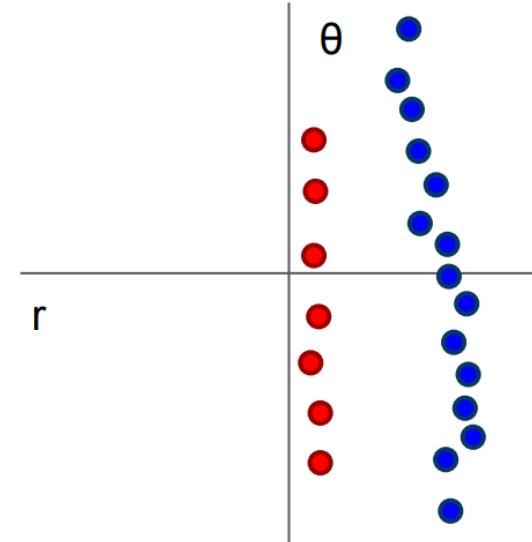
Does it generalize???

A: Yes, by transforming the data coordinates *before* classification

Image Features: Motivation



$$f(x, y) = (r(x, y), \theta(x, y))$$

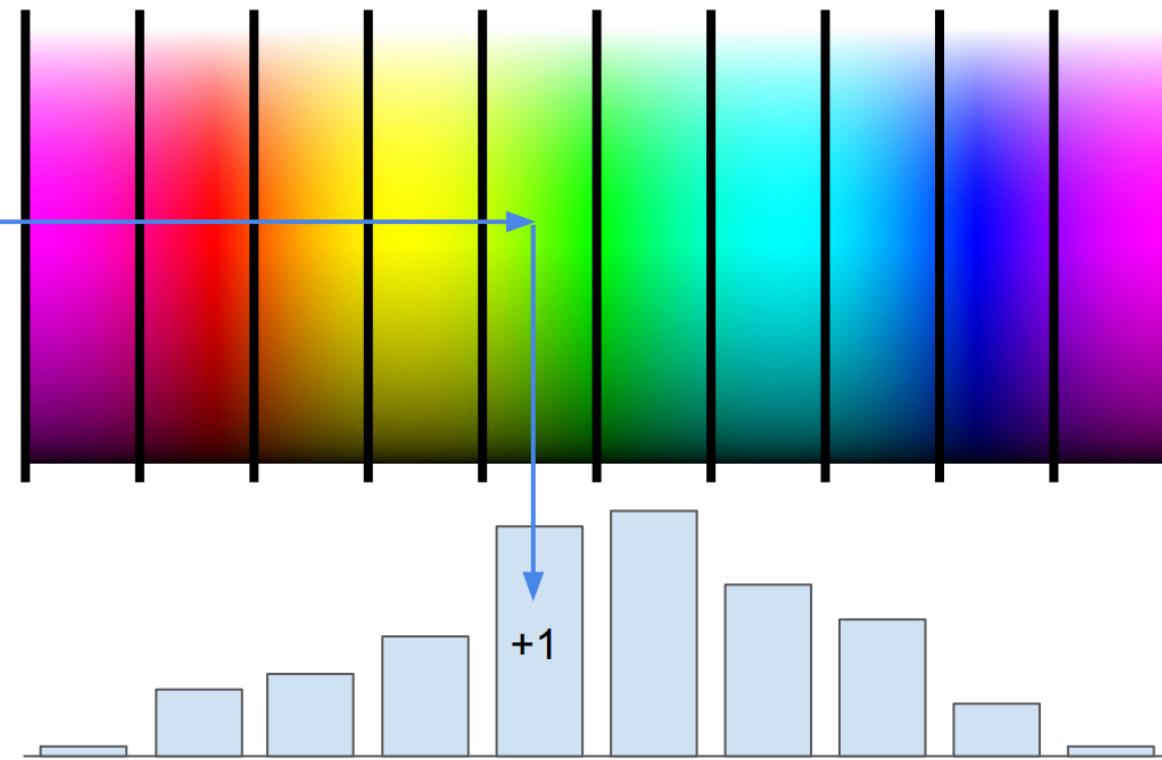


Cannot separate red and blue points with linear classifier

After applying feature transform, points can be separated by linear classifier

From Stanford CS231: <http://cs231n.stanford.edu/>

Example: Color Histogram



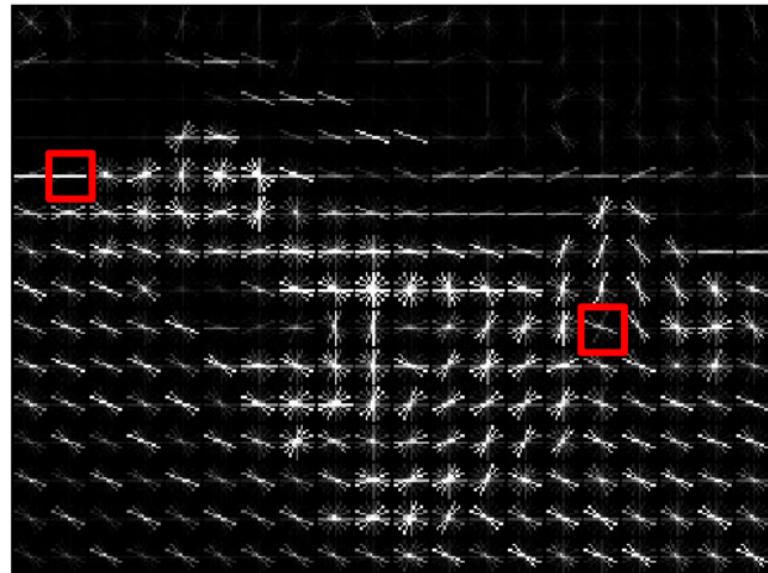
From Stanford CS231: <http://cs231n.stanford.edu/>

Example: Histogram of Oriented Gradients (HoG)



Divide image into 8x8 pixel regions
Within each region quantize edge
direction into 9 bins

Lowe, "Object recognition from local scale-invariant features", ICCV 1999
Dalal and Triggs, "Histograms of oriented gradients for human detection," CVPR 2005



Example: 320x240 image gets divided
into 40x30 bins; in each bin there are
9 numbers so feature vector has
 $30 \times 40 \times 9 = 10,800$ numbers

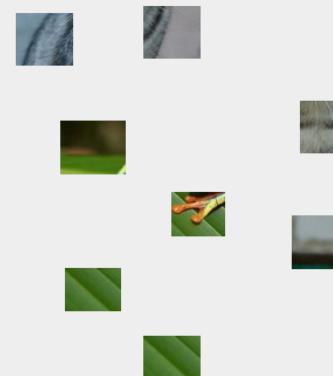
From Stanford CS231: <http://cs231n.stanford.edu/>

Example: Bag of Words

Step 1: Build codebook



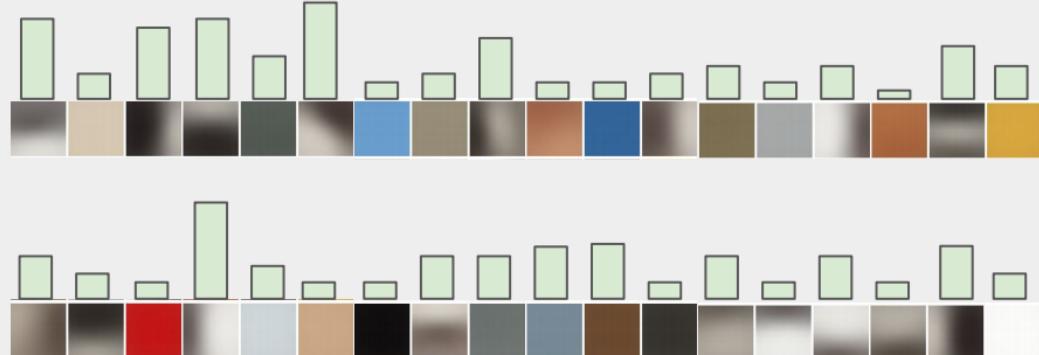
Extract random patches



Cluster patches to
form “codebook”
of “visual words”



Step 2: Encode images



Fei-Fei and Perona, “A bayesian hierarchical model for learning natural scene categories”, CVPR 2005

From Stanford CS231: <http://cs231n.stanford.edu/>

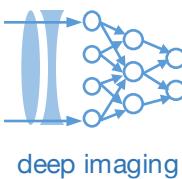
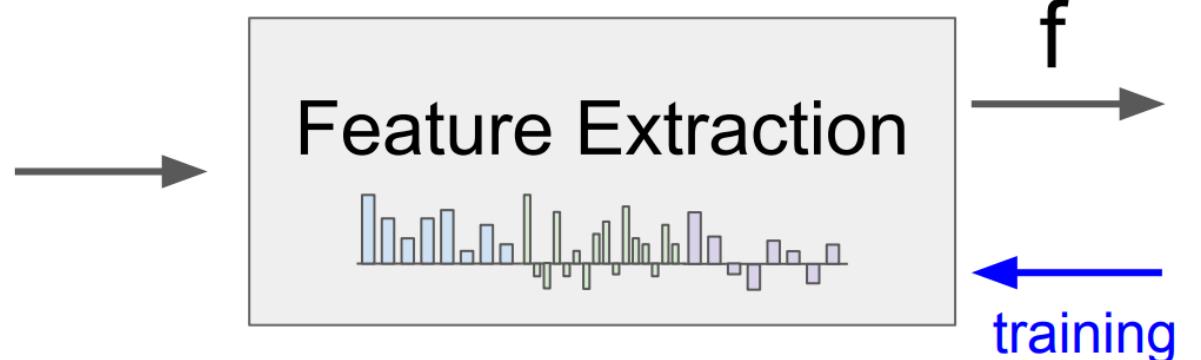
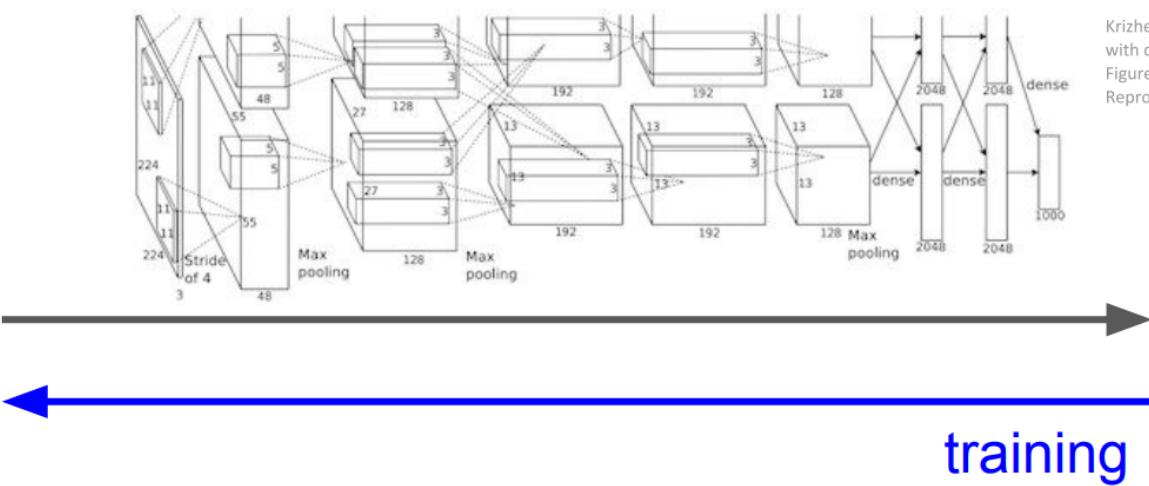


Image features vs ConvNets

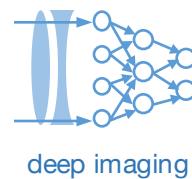


10 numbers giving scores for classes



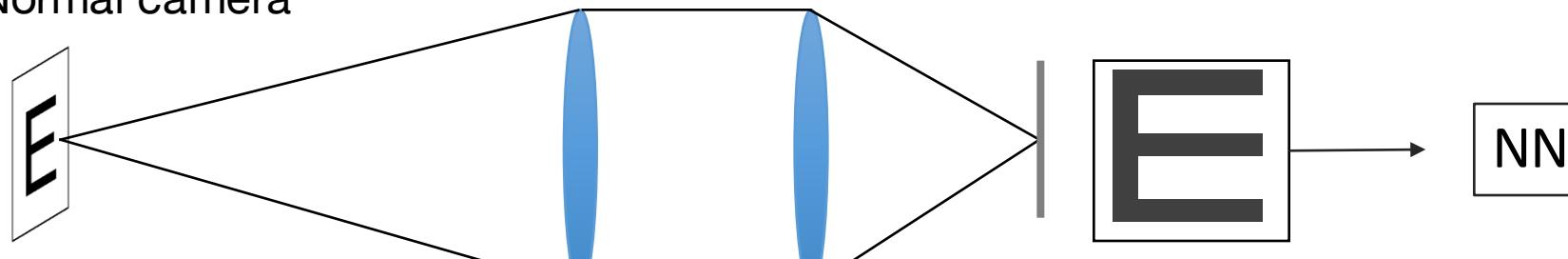
10 numbers giving scores for classes

From Stanford CS231: <http://cs231n.stanford.edu/>



Hand-crafted versus learned features also applies to imaging

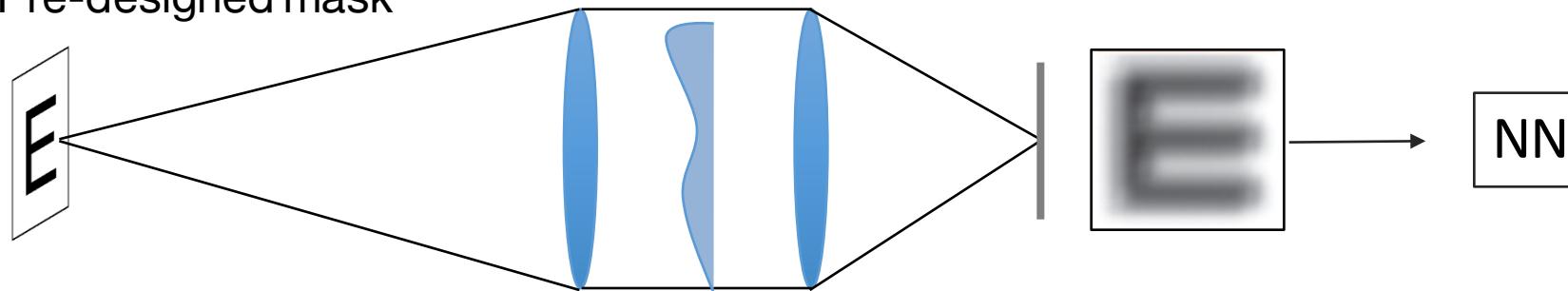
Normal camera



Classification acc.:

NN → 80%

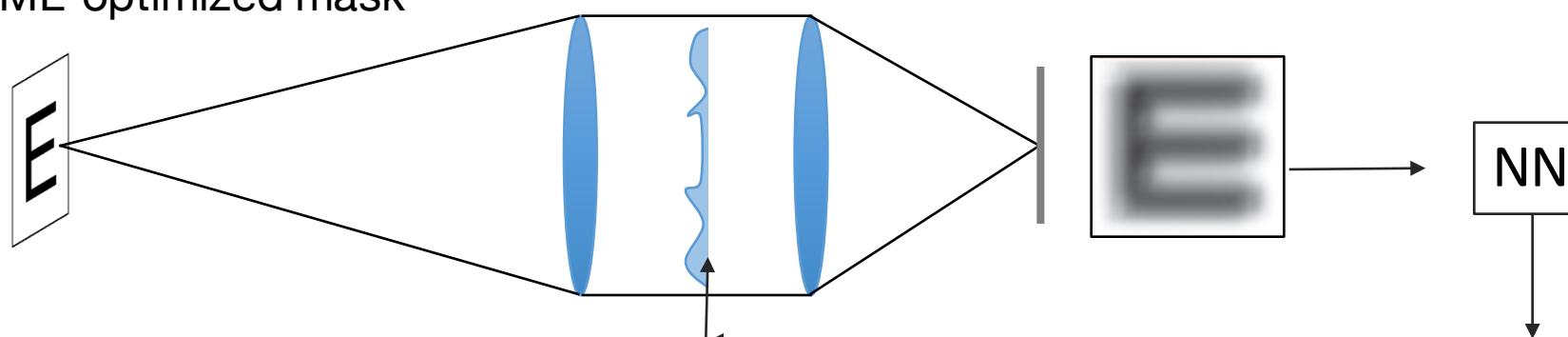
Pre-designed mask



Classification acc.:

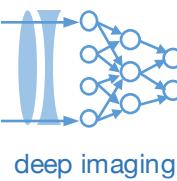
NN → 90%

ML-optimized mask



Classification acc.:

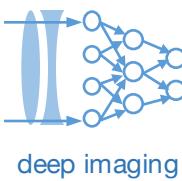
NN → 95%



Statistical Machine Learning in 30 minutes

Two competing goals in machine learning:

1. Can we make sure the in-sample error $L_{in}(y, f(x, W))$ is small enough?
 - Appropriate cost function
 - “complex enough” model



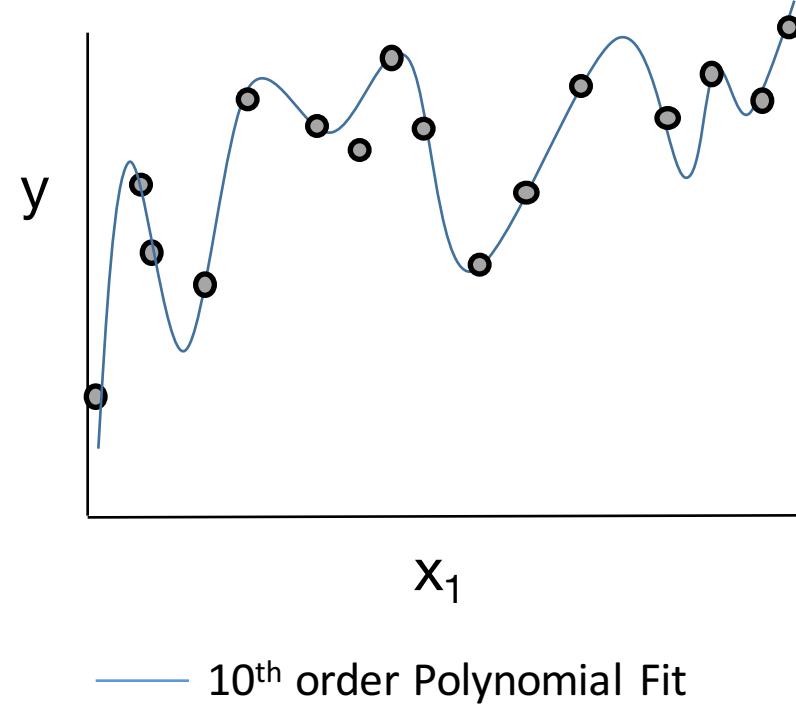
Statistical Machine Learning in 30 minutes

Two competing goals in machine learning:

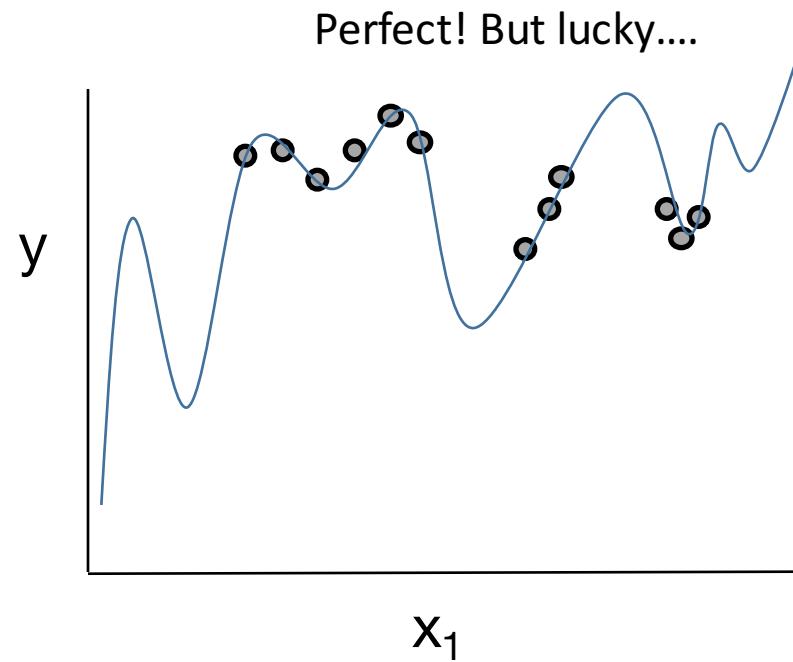
1. Can we make sure the in-sample error $L_{in}(y, f(x, W))$ is small enough?
 - Appropriate cost function
 - “complex enough” model
2. Can we make sure that $L_{out}(y, f(x, W))$ is close enough to $L_{in}(y, f(x, W))$?
 - Probabilistic analysis says yes!
 - $|L_{in} - L_{out}|$ bounded from above
 - Bound grows with model capacity (bad)
 - Bound shrinks with # of training examples (good)

Model overfitting versus underfitting – a thought exercise

Let's fit these “training” data points:

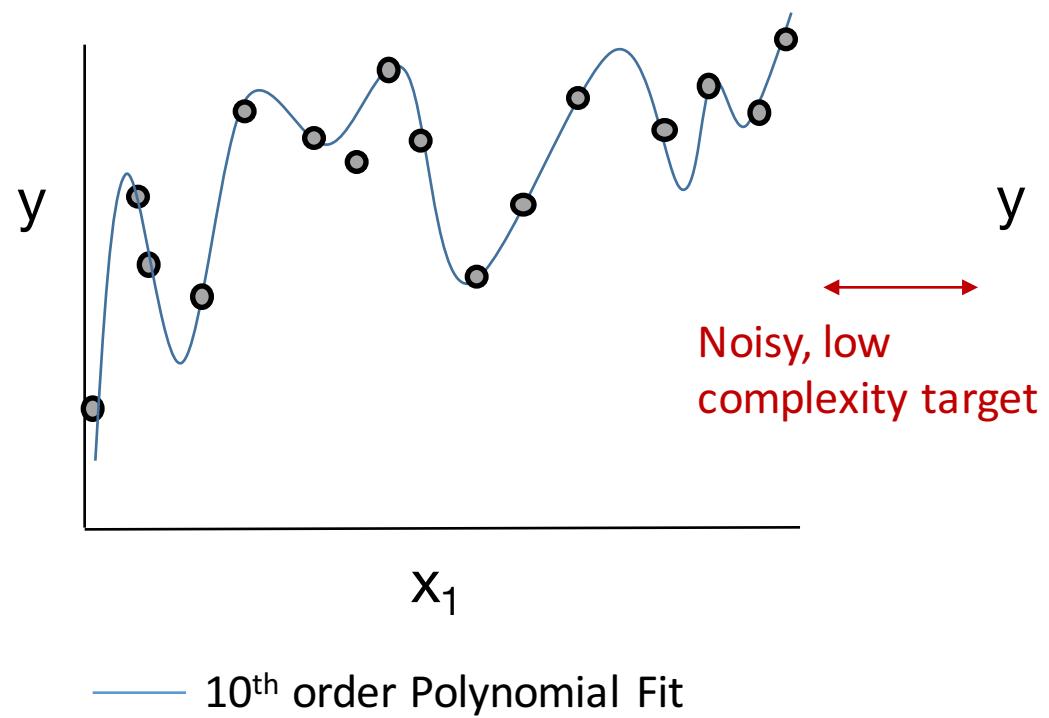


And then here's our testing dataset – good?

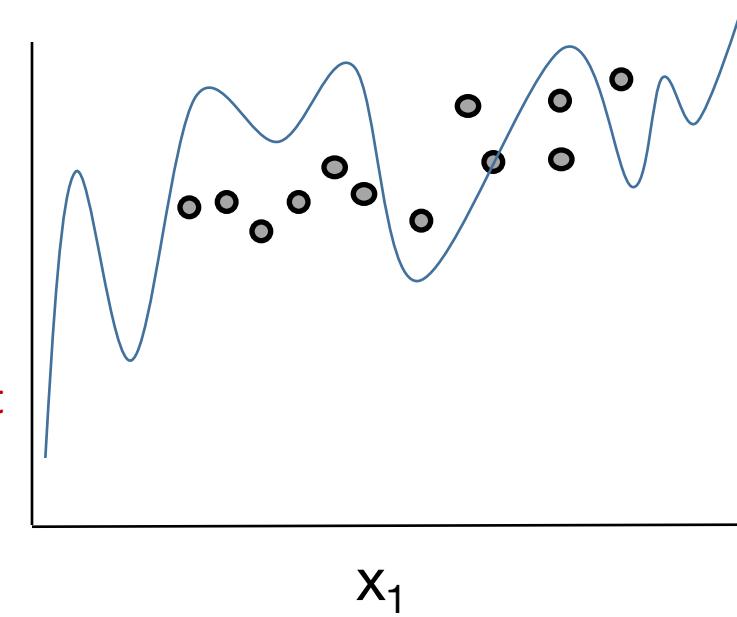


Model overfitting versus underfitting – a thought exercise

Let's fit these "training" data points:

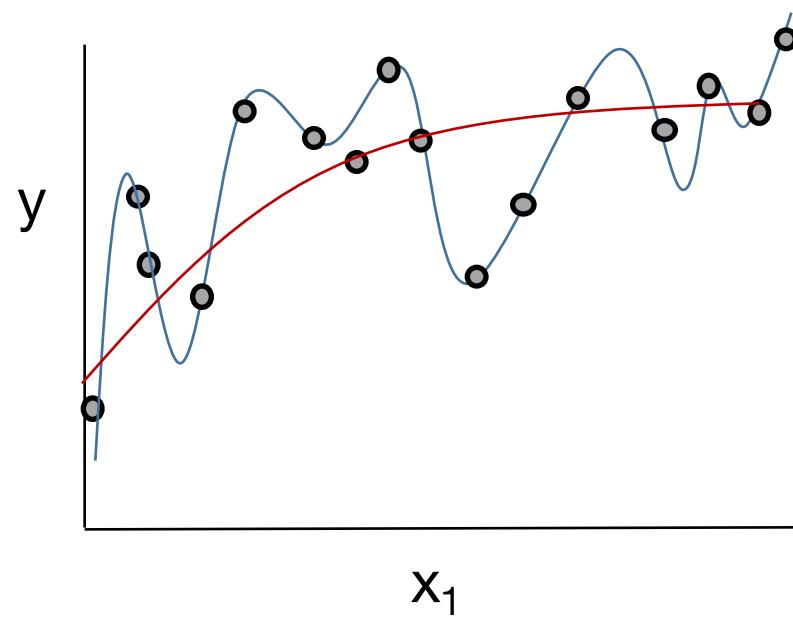


What if our test dataset was this :



Model overfitting versus underfitting – a thought exercise

Let's fit these "training" data points:



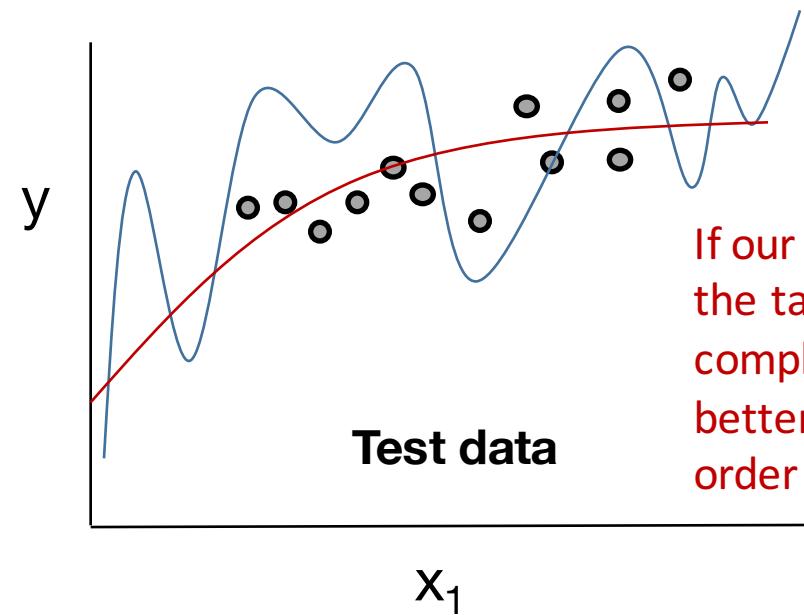
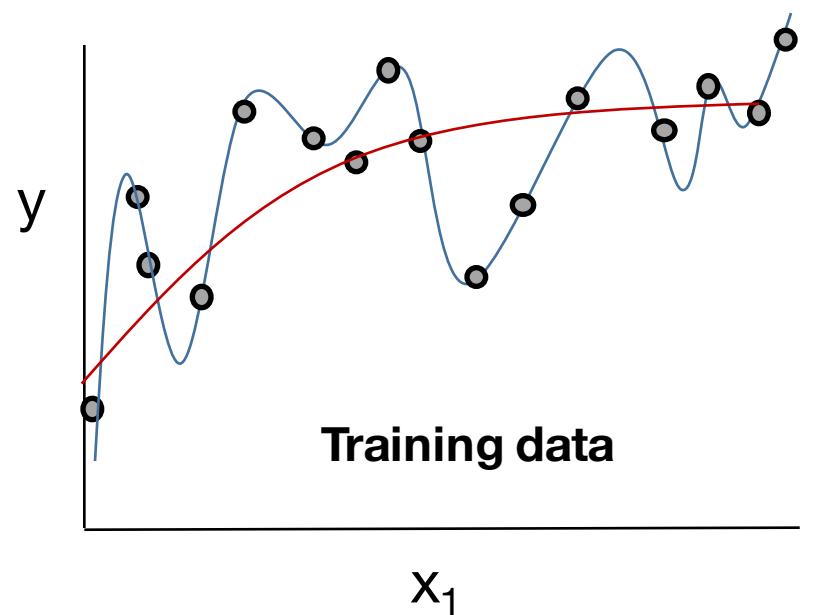
- 10th order Polynomial Fit
- 2nd order Polynomial Fit

What if our test dataset was this :



If our data was noisy and the target followed a low-complexity model, we'd be better off with a second order fit!

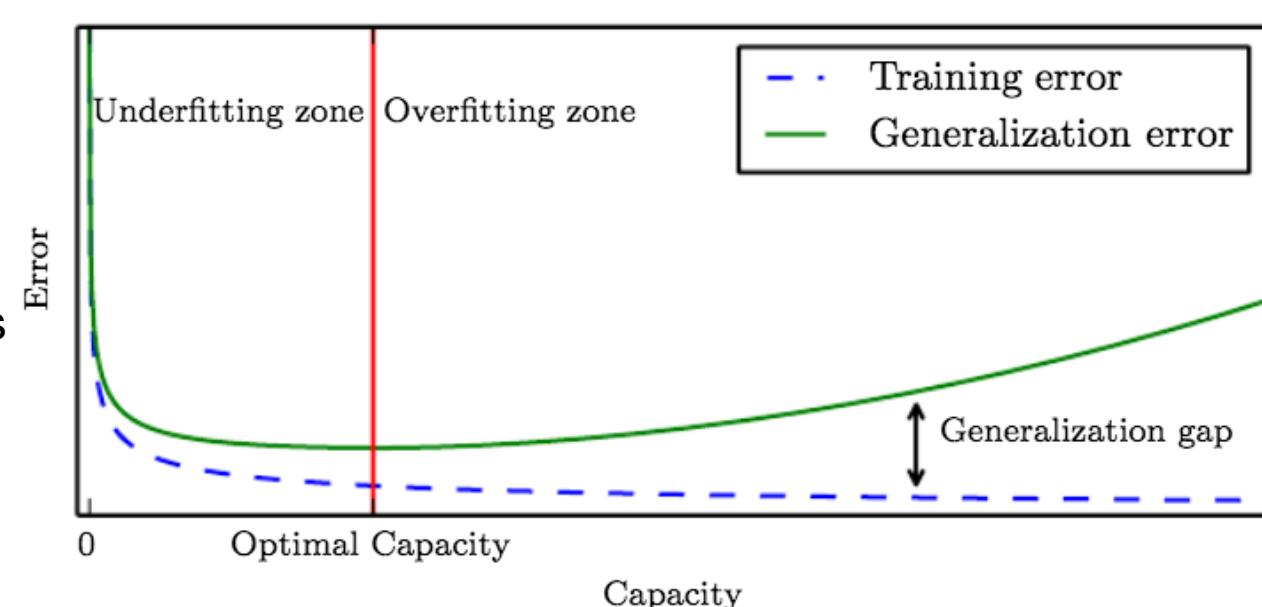
Model overfitting versus underfitting – a thought exercise



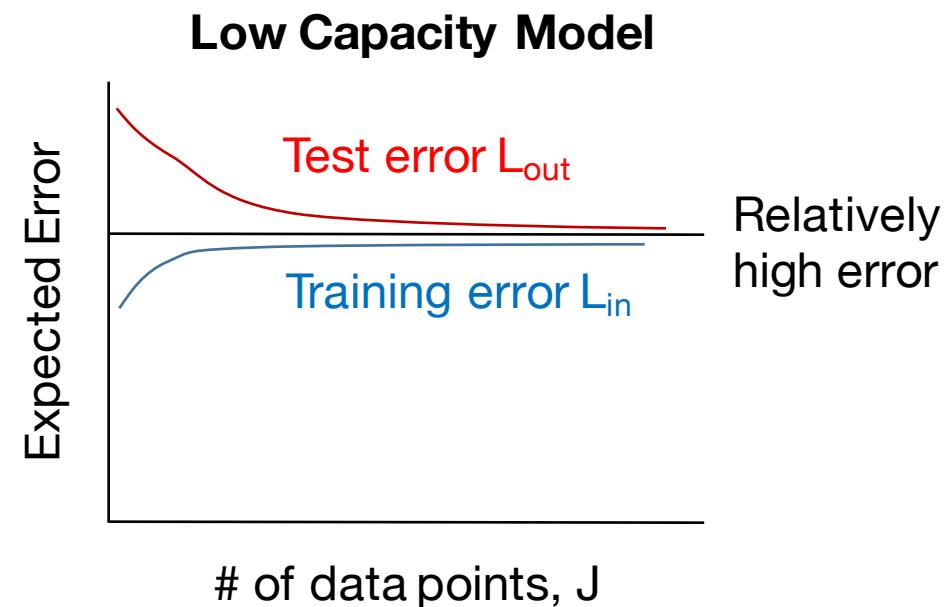
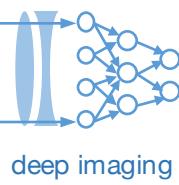
Model capacity: ability to fit a wide range of functions

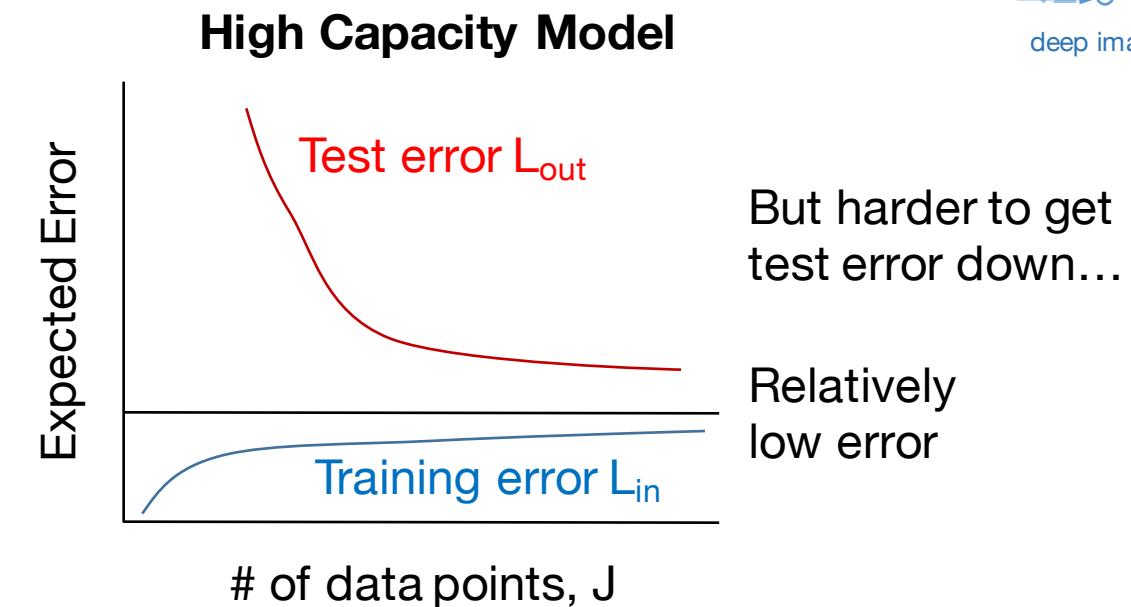
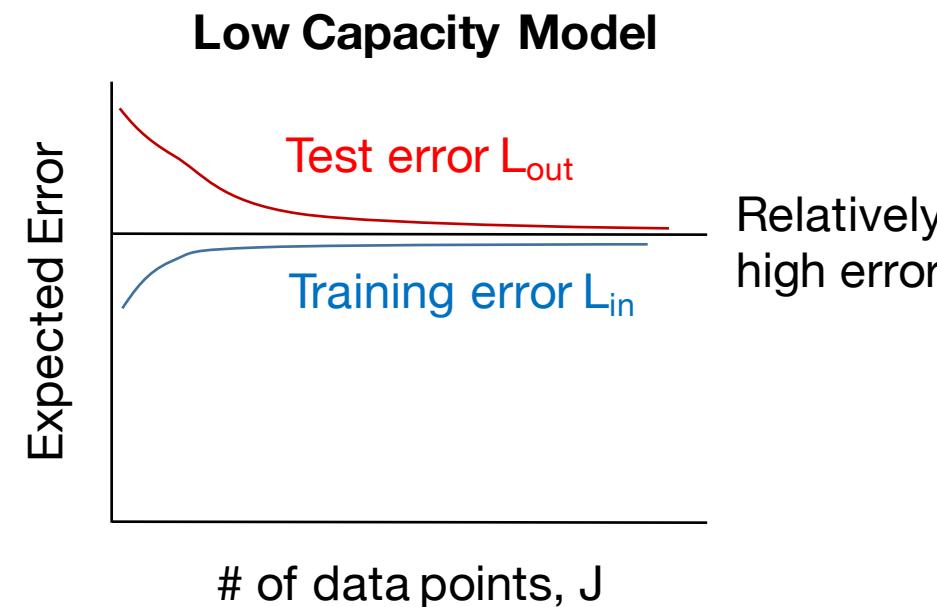
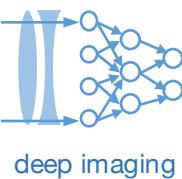
Control capacity through model's hypothesis space (set of functions model can take)

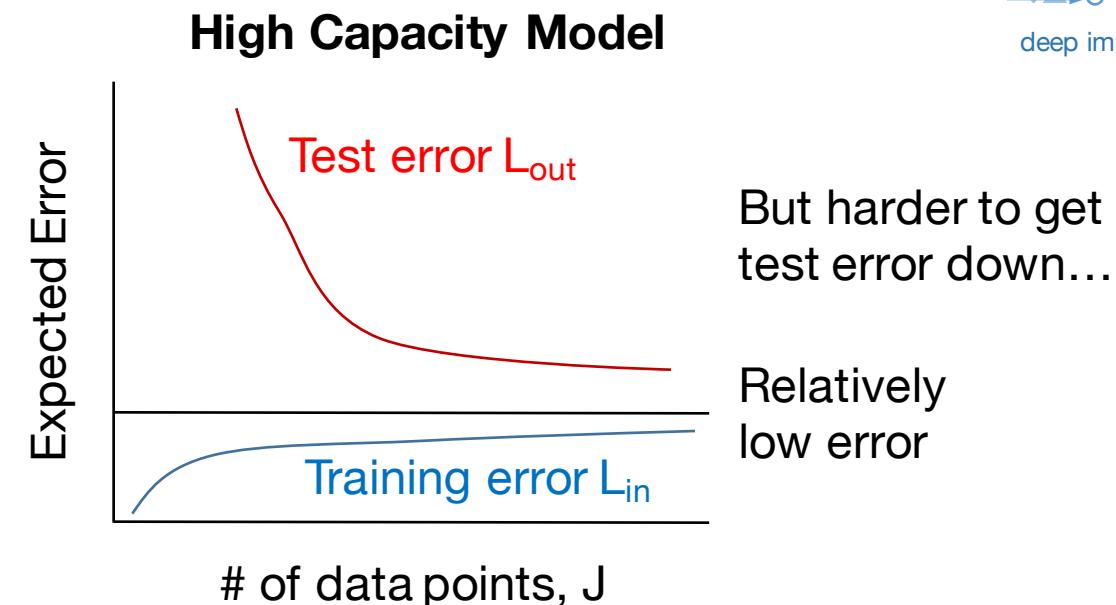
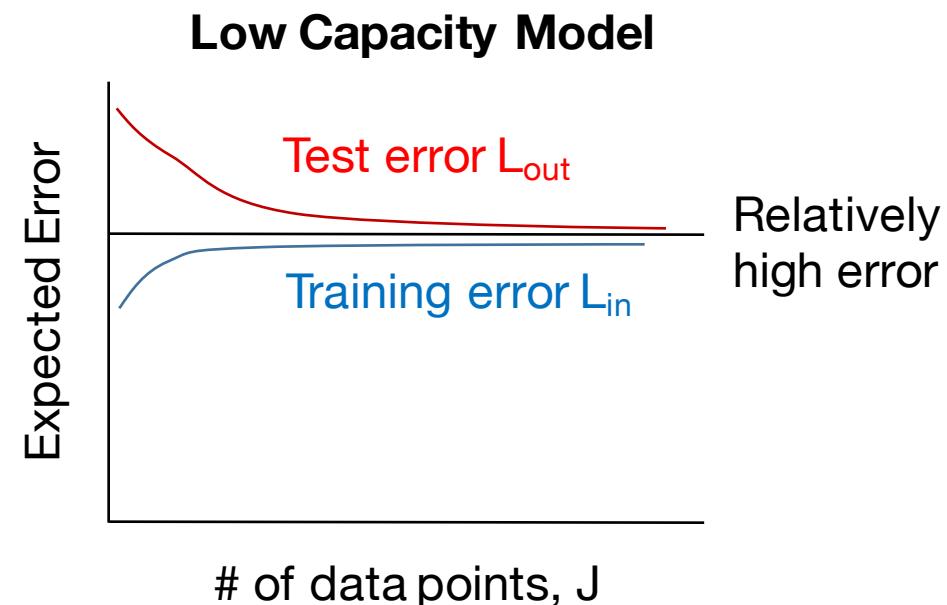
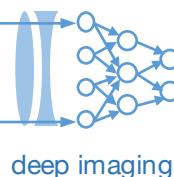
Hard to know ahead of time!



Deep Learning, I. Goodfellow et al., Fig. 5.3

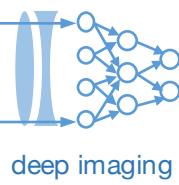




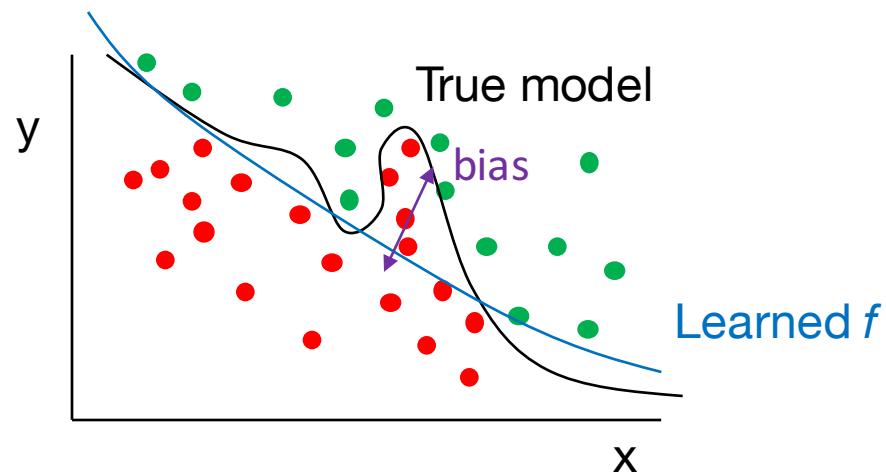
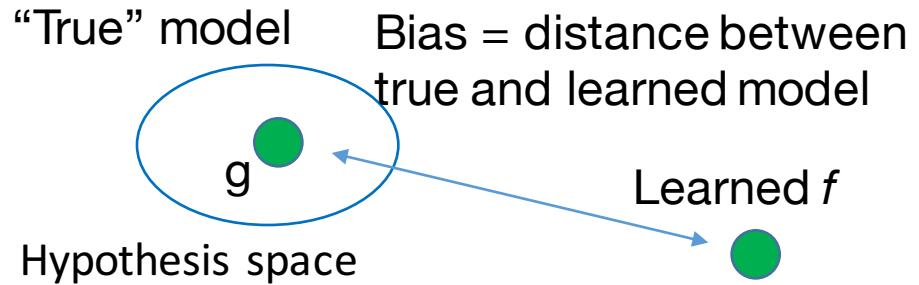


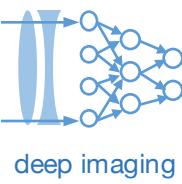
Take away concepts:

- Can't ever really expect test error to be less than training error
- Complicated models tend to appear to “do better” during training, before trying test data
- When the model gets complicated and you don’t have enough data, challenging to get test error down

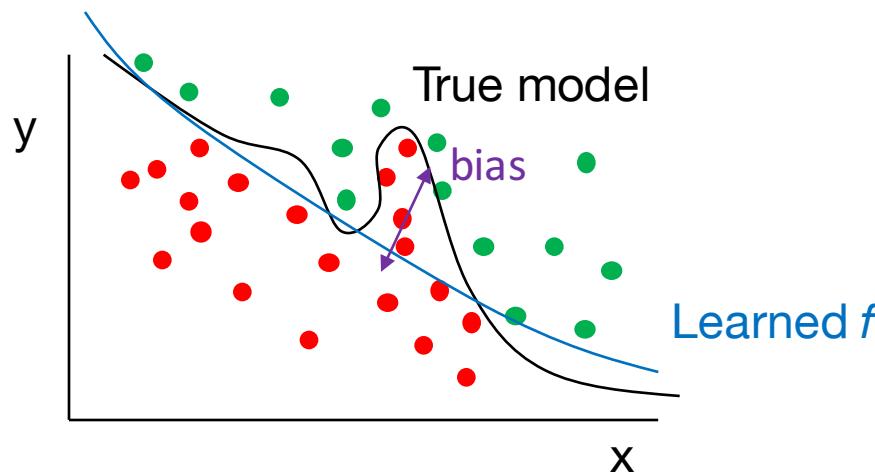
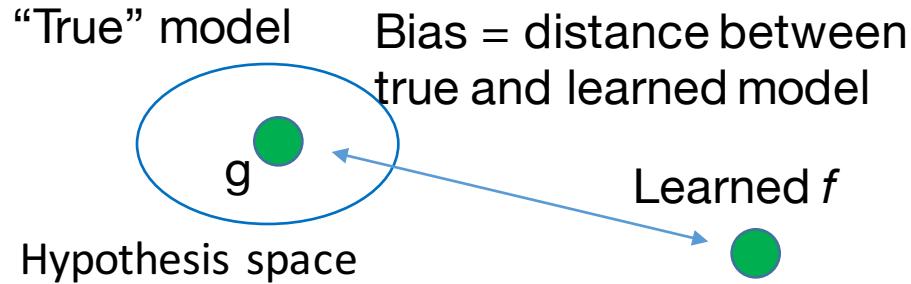


Model bias versus variance





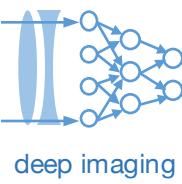
Model bias versus variance



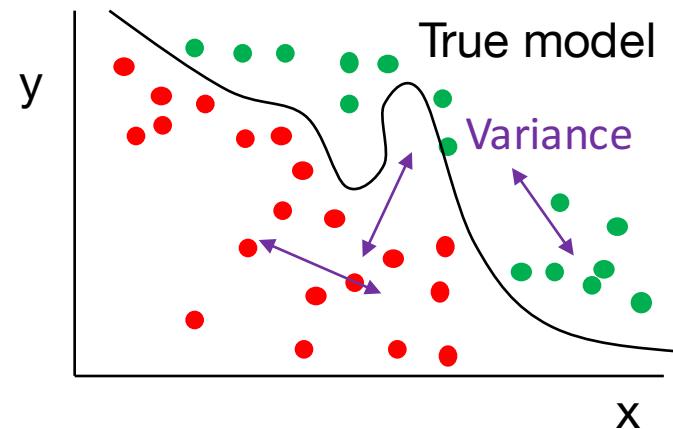
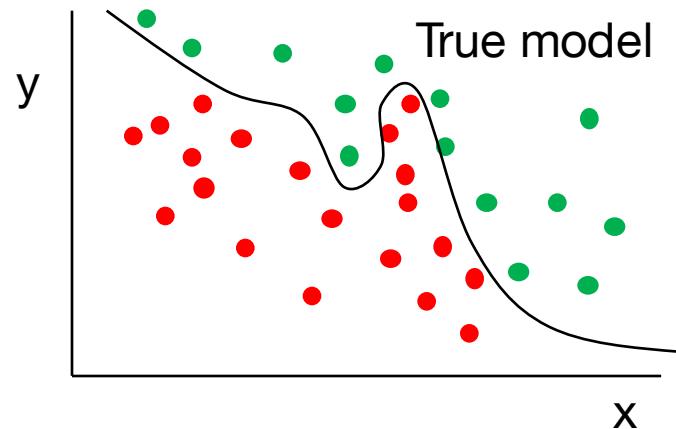
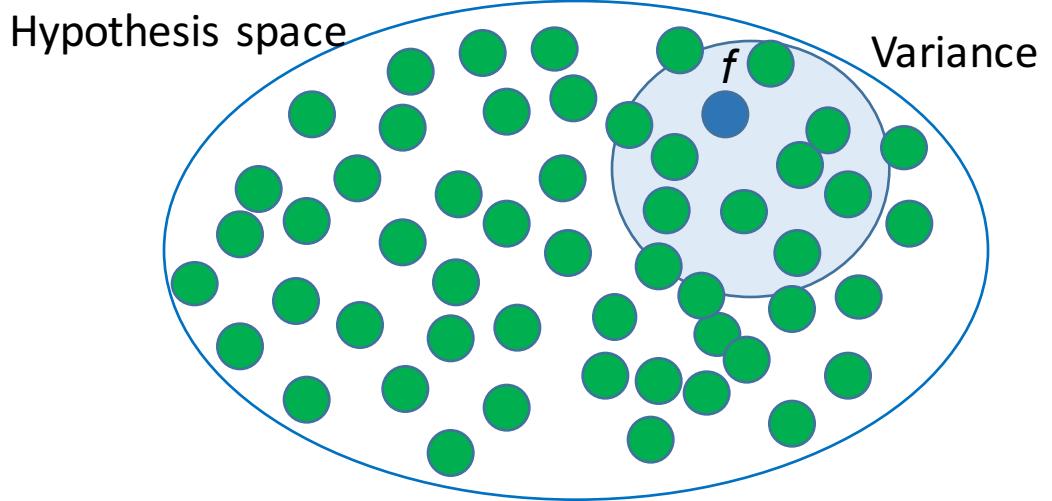
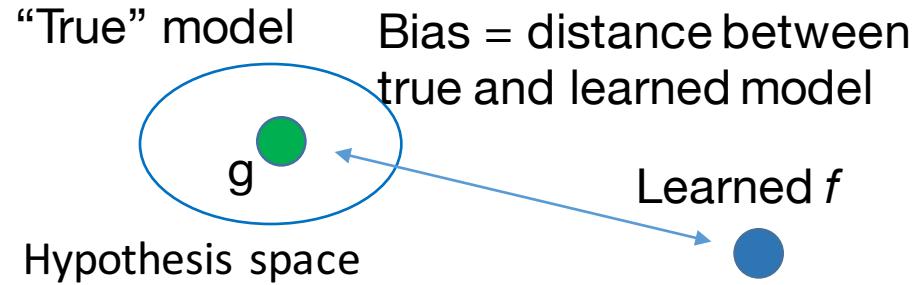
$$\text{Bias} = (g(x) - f(x))^2$$

Measures how far our learning model f is biased away from target function g (for perfect training data classification)

Models that tend to be “a bit too simple” are biased away from “true” model

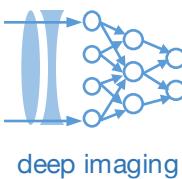


Model bias versus variance

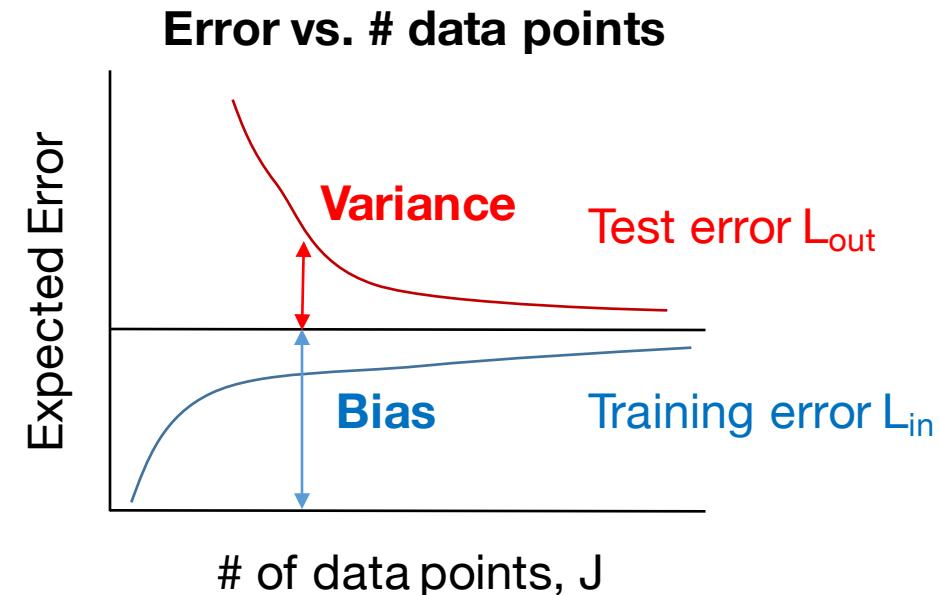
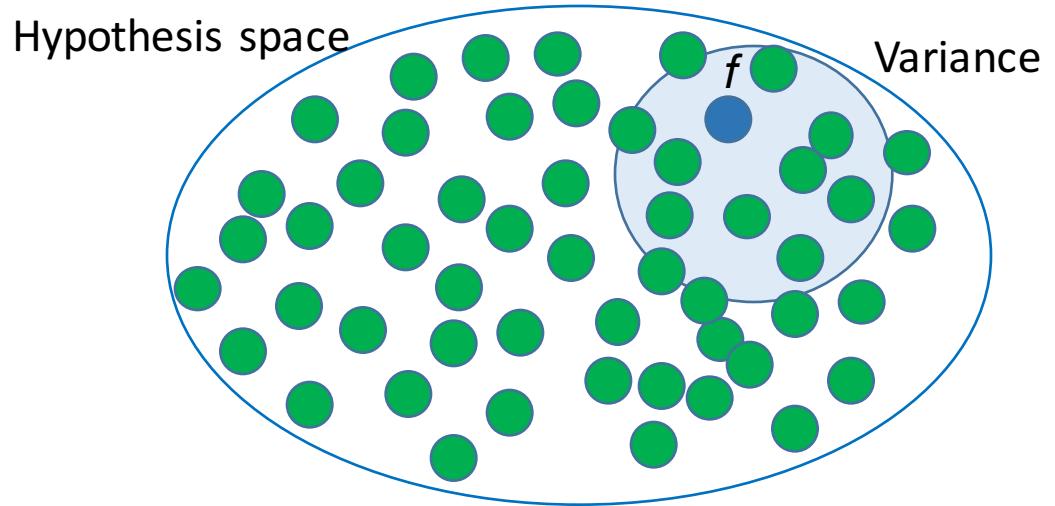
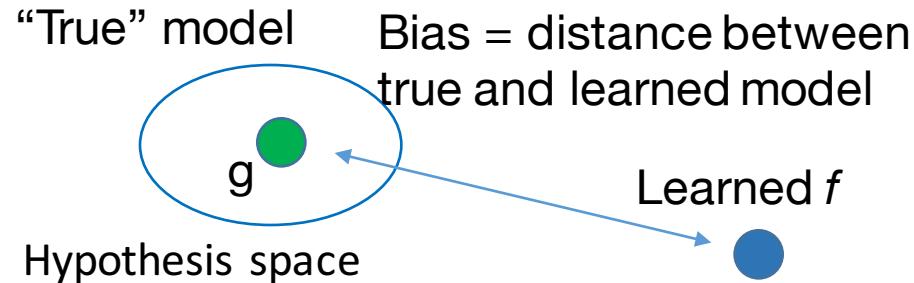


$$\text{Variance} = \text{Var}[g(x)]$$

More complicated datasets exhibit lots of variance between training and test set

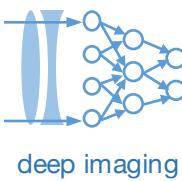


Model bias versus variance



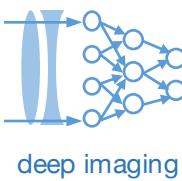
Test Error is sum of model bias and variance!

Goal is to find a model f that balances between these two quantities for a given dataset



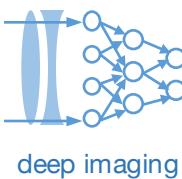
How to formally define capacity and complexity?

- Short answer: it's complicated...
- Related to something called the *VC Dimension*
 - Can provide theoretical bounds on performance
 - Dimensional bounds rather than scalar bounds...
- I decided not to go into it, but please let me know if you'd like me to!



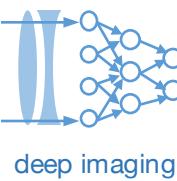
Conclusions from statistical machine learning

- Conclusion: you want a model that is complex enough to capture variations within high-dimensional space, but not too complex such that it overfits the data
- Want a model with a high capacity, but can still *generalize* to data outside training set
 - More data -> less overfitting, complex target -> more overfitting
- For simple models, we can measure complexity via degrees of freedom, the VC bound and so-on to help us nail down ideal models that can generalize well



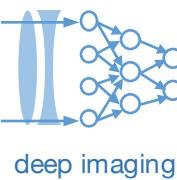
Conclusions from statistical machine learning

- Conclusion: you want a model that is complex enough to capture variations within high-dimensional space, but not too complex such that it overfits the data
- Want a model with a high capacity, but can still *generalize* to data outside training set
 - More data -> less overfitting, complex target -> more overfitting
- For simple models, we can measure complexity via degrees of freedom, the VC bound and so-on to help us nail down ideal models that can generalize well
- **For DL models:** this will get too hard...here's a few counter-intuitive properties:
 1. A fixed DL *architecture* exhibits data-dependent complexities
 - e.g., “good” DL networks achieve 0 training error on images with random labels, so cannot generalize at all in this case, and are too complex
 2. DL networks with more hidden units leads to *better* generalization (the main finding of the last few years). So deeper models tend to be less complex, actually...
 3. Complexity depends upon loss function and optimization method...



Important to remember: “No Free Lunch Theorem”

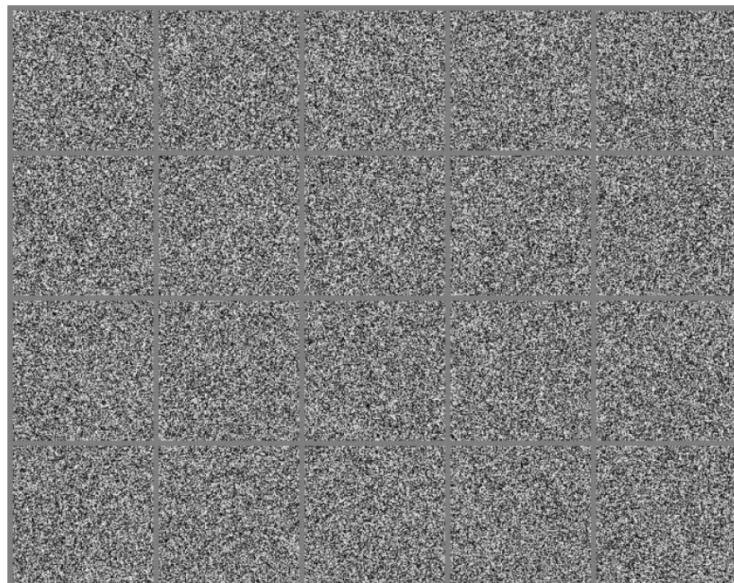
- “Averaged over all possible data-generating distributions, every classification algorithm has the same error rate when classifying previously unobserved points.”
- The most sophisticated DL algorithm has save average performance (averaged over all possible tasks) as the simplest.



Important to remember: “No Free Lunch Theorem”

- “Averaged over all possible data-generating distributions, every classification algorithm has the same error rate when classifying previously unobserved points.”
- The most sophisticated DL algorithm has save average performance (averaged over all possible tasks) as the simplest.
- Must make assumptions about probability distributions of inputs we’ll encounter in real-world

Set of 20 “images”, random Gaussian distribution



Face at different orientations =
manifold n-D space

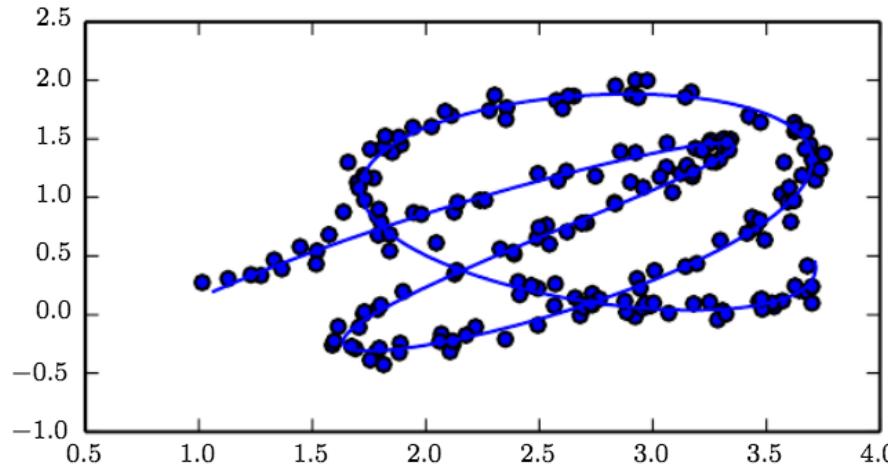


Deep Learning, I. Goodfellow et al., Fig. 5.12-13

Important to remember: “No Free Lunch Theorem”

- “Averaged over all possible data-generating distributions, every classification algorithm has the same error rate when classifying previously unobserved points.”
- The most sophisticated DL algorithm has save average performance (averaged over all possible tasks) as the simplest.
- Must make assumptions about probability distributions of inputs we’ll encounter in real-world

1D Manifold in 2D space

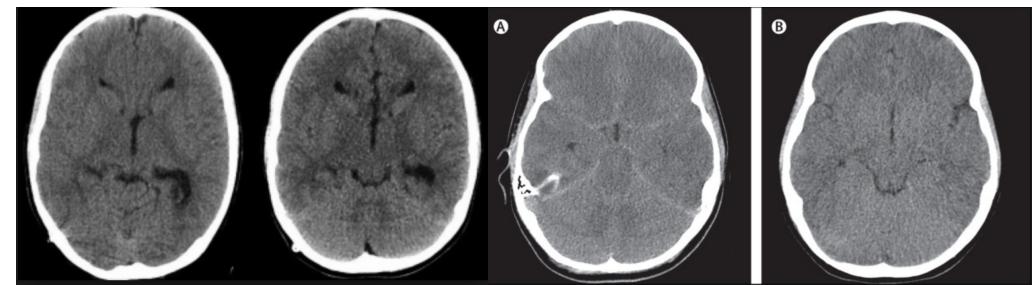


Deep Learning, I. Goodfellow et al., Fig. 5.11

Manifold
Hypothesis



CT reconstructions of every brain in the world = kD manifold in nD space?



...