

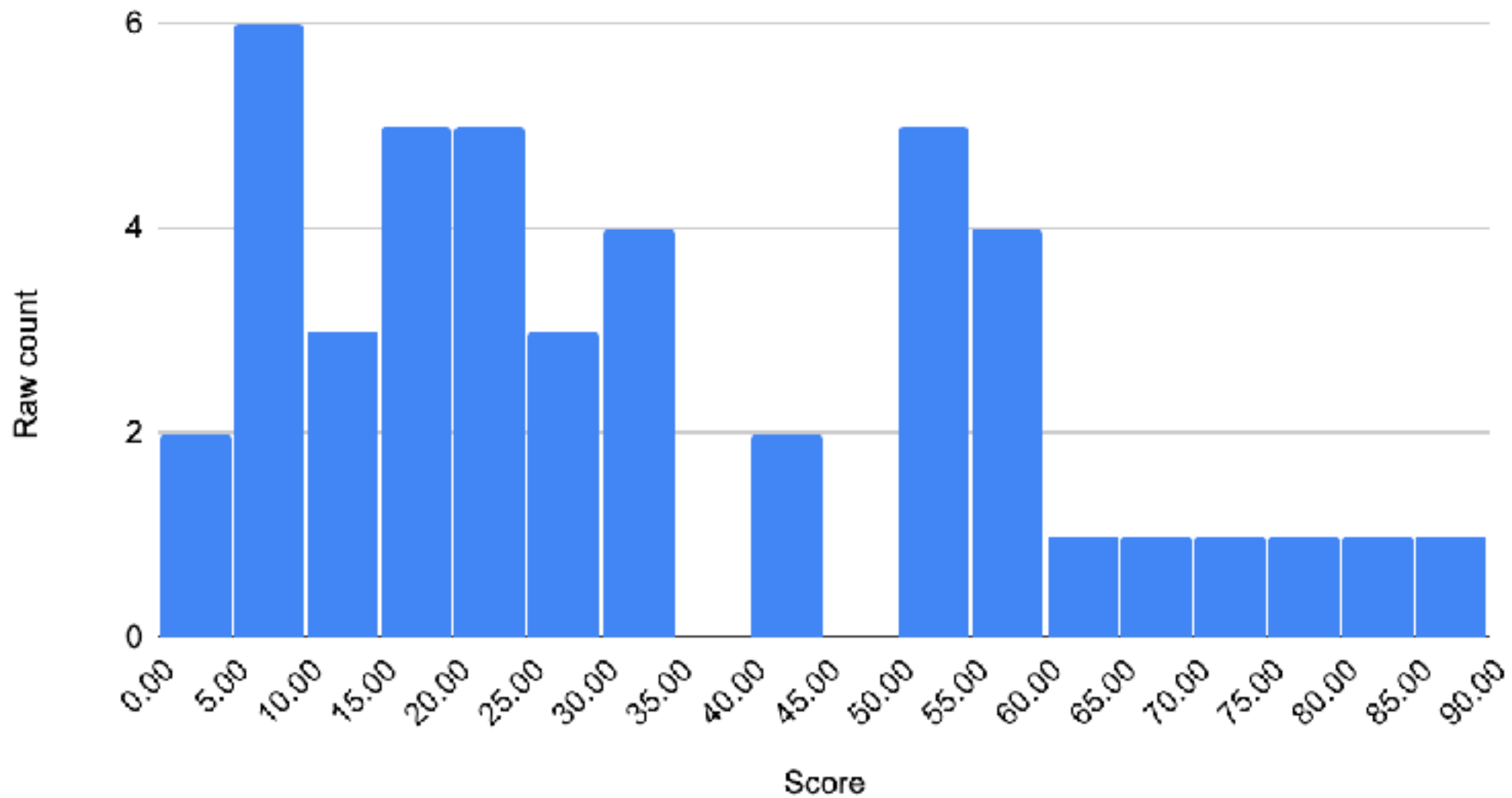
95-865 Unstructured Data Analytics

Week 5: Intro to predictive data analytics, neural nets, and deep learning

George Chen

Quiz 1

Fall 2019 95-865 Quiz 1 Histogram



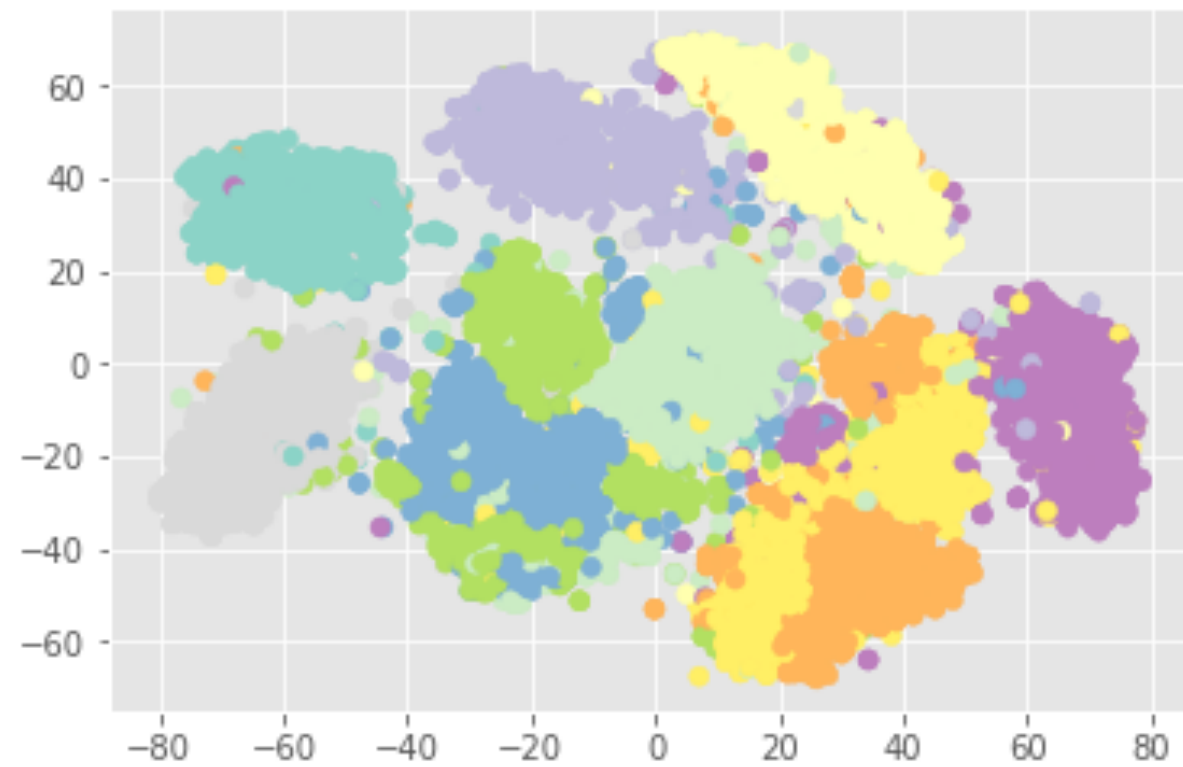
Mean: 33.3, std dev: 23.6, max achieved: 87

Quiz 1 Regrade Requests

- How regrades work:
 1. Study solutions (already posted in Canvas under “Files”) very carefully
 2. If you think there’s a mistake, send me an email and be very specific about what was incorrectly graded and how many points are at stake
 3. We will regrade your whole quiz 1 (the version that you submitted to Canvas on the quiz day)
 4. Your score can go up, go down, or stay the same, and the regraded result is final
- Due this Friday 11:59pm Pittsburgh time

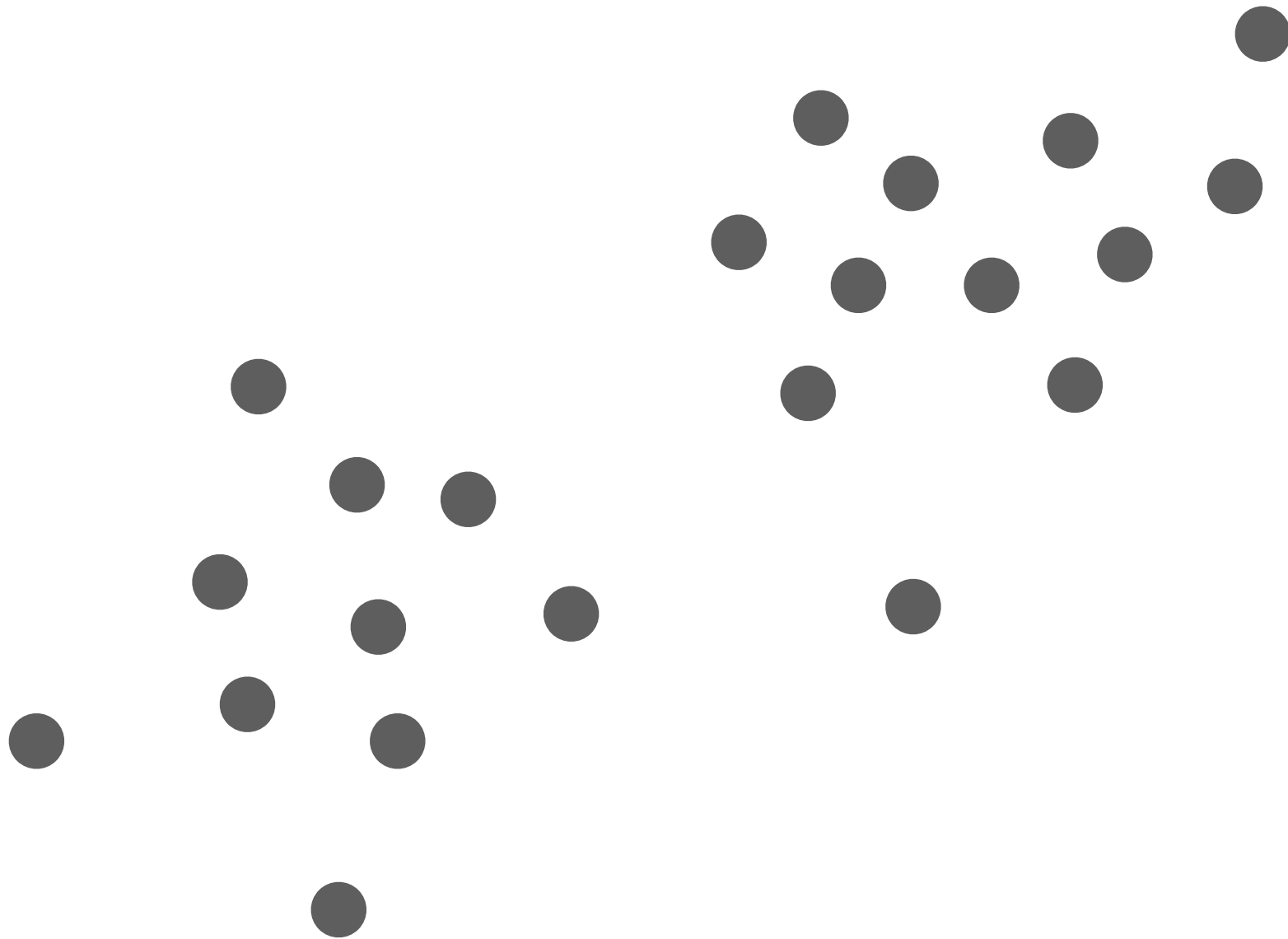
Disclaimer: unfortunately “*k*”
means many things

What if we have labels?



Example: MNIST handwritten digits have known labels

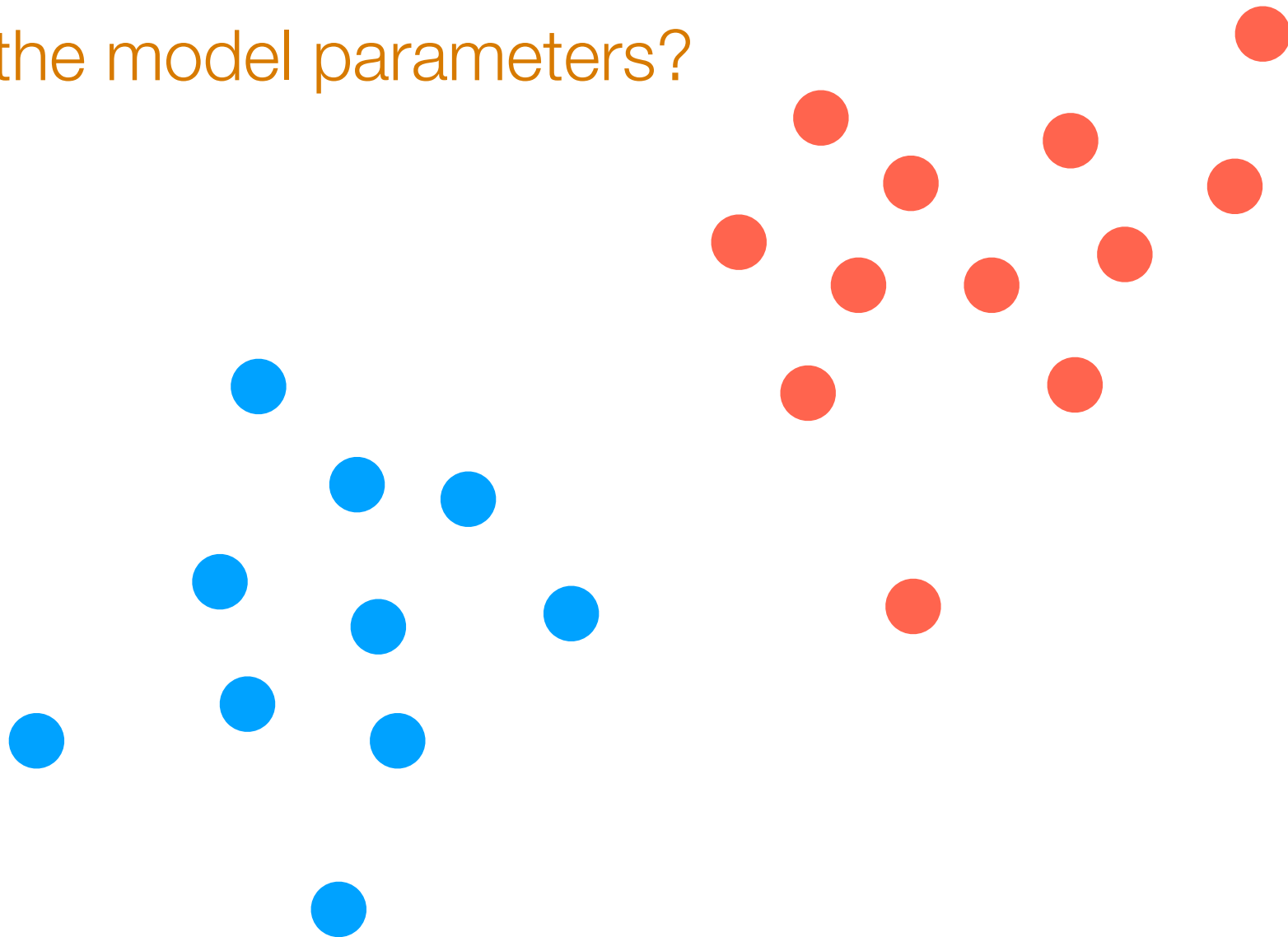
If the labels are known...



If the labels are known...

And we assume data generated by GMM...

What are the model parameters?



Flashback: Learning a GMM

Don't need this top part if we know the labels!

Step 0: Pick k

Step 1: Pick guesses for **cluster probabilities, means, and covariances** (often done using k -means)

Repeat until convergence:

Step 2: Compute probability of each point belonging to each of the k clusters

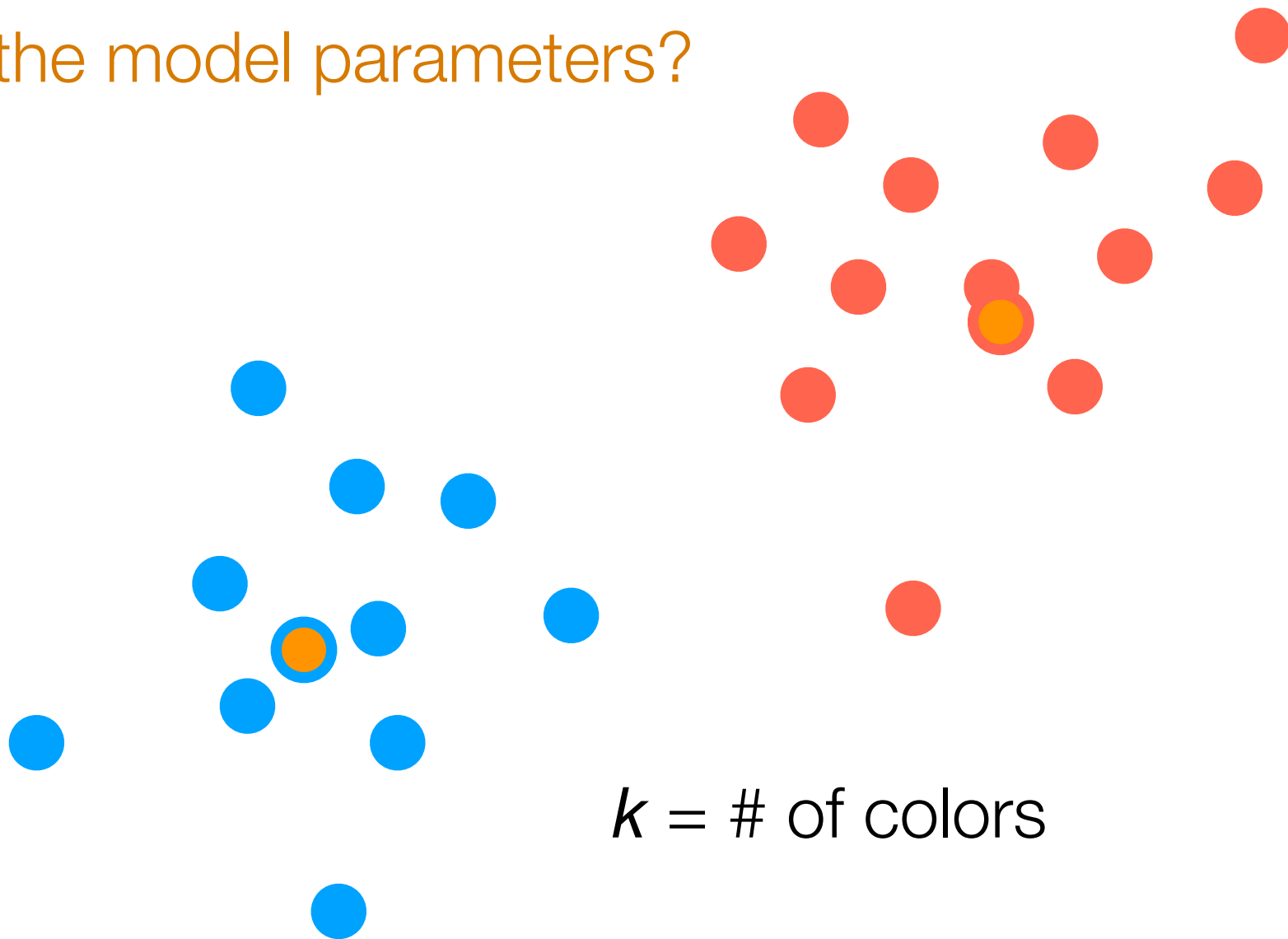
Step 3: Update **cluster probabilities, means, and covariances** carefully accounting for probabilities of each point belonging to each of the clusters

We don't need to repeat until convergence

If the labels are known...

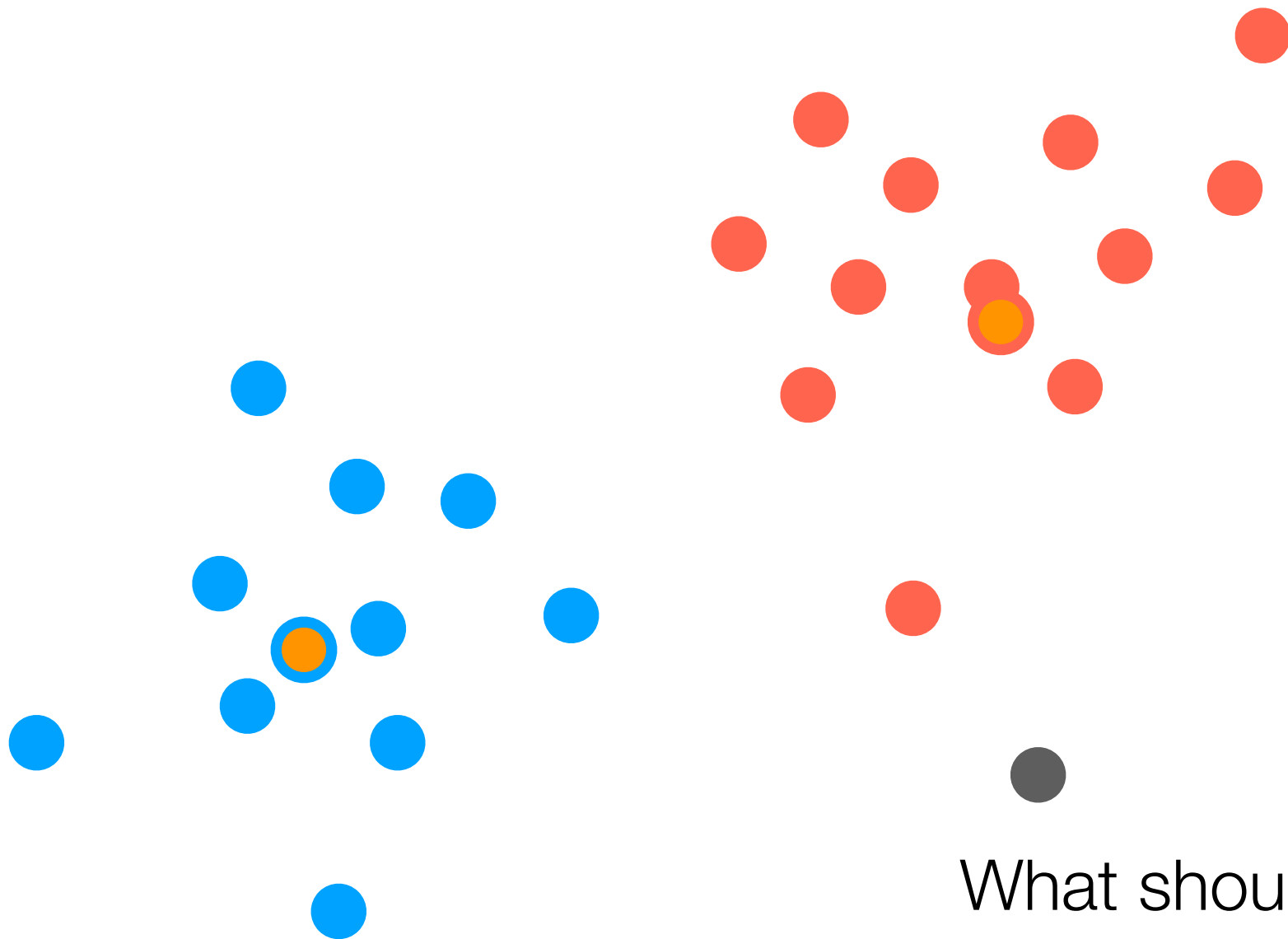
And we assume data generated by GMM...

What are the model parameters?



$k = \#$ of colors

We can directly estimate
cluster means, covariances

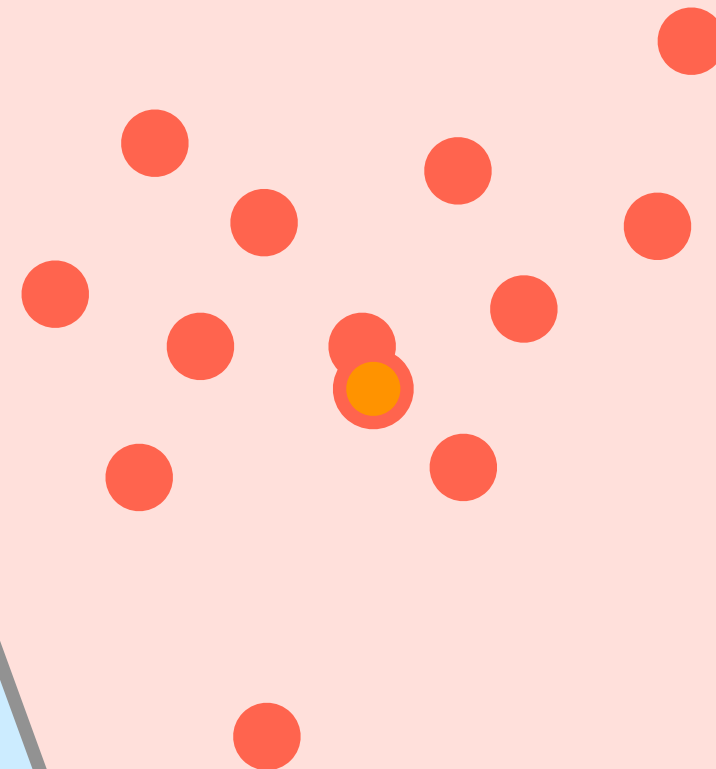
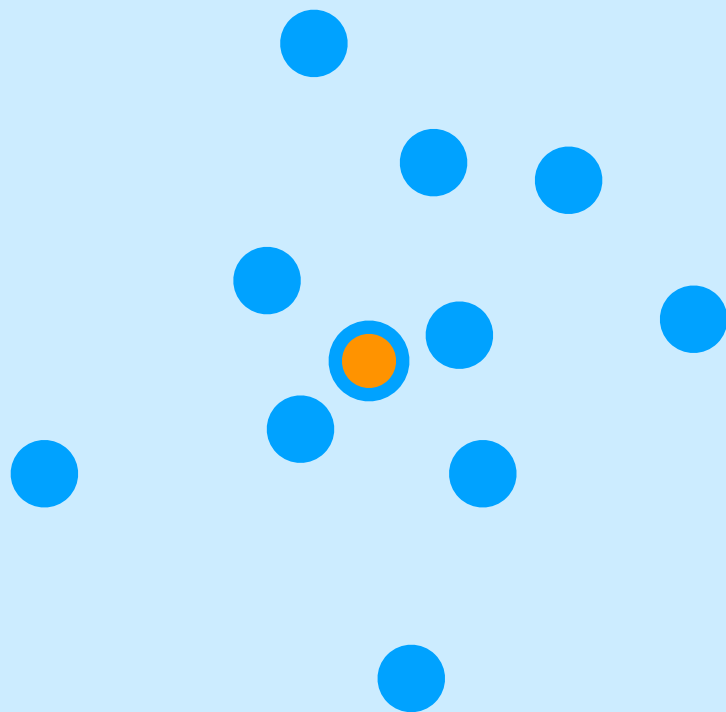


What should the label of
this new point be?

Whichever cluster has
higher probability!

We just created a **classifier**
(a procedure that given a new data point tells us what “class” it belongs to)

Decision boundary



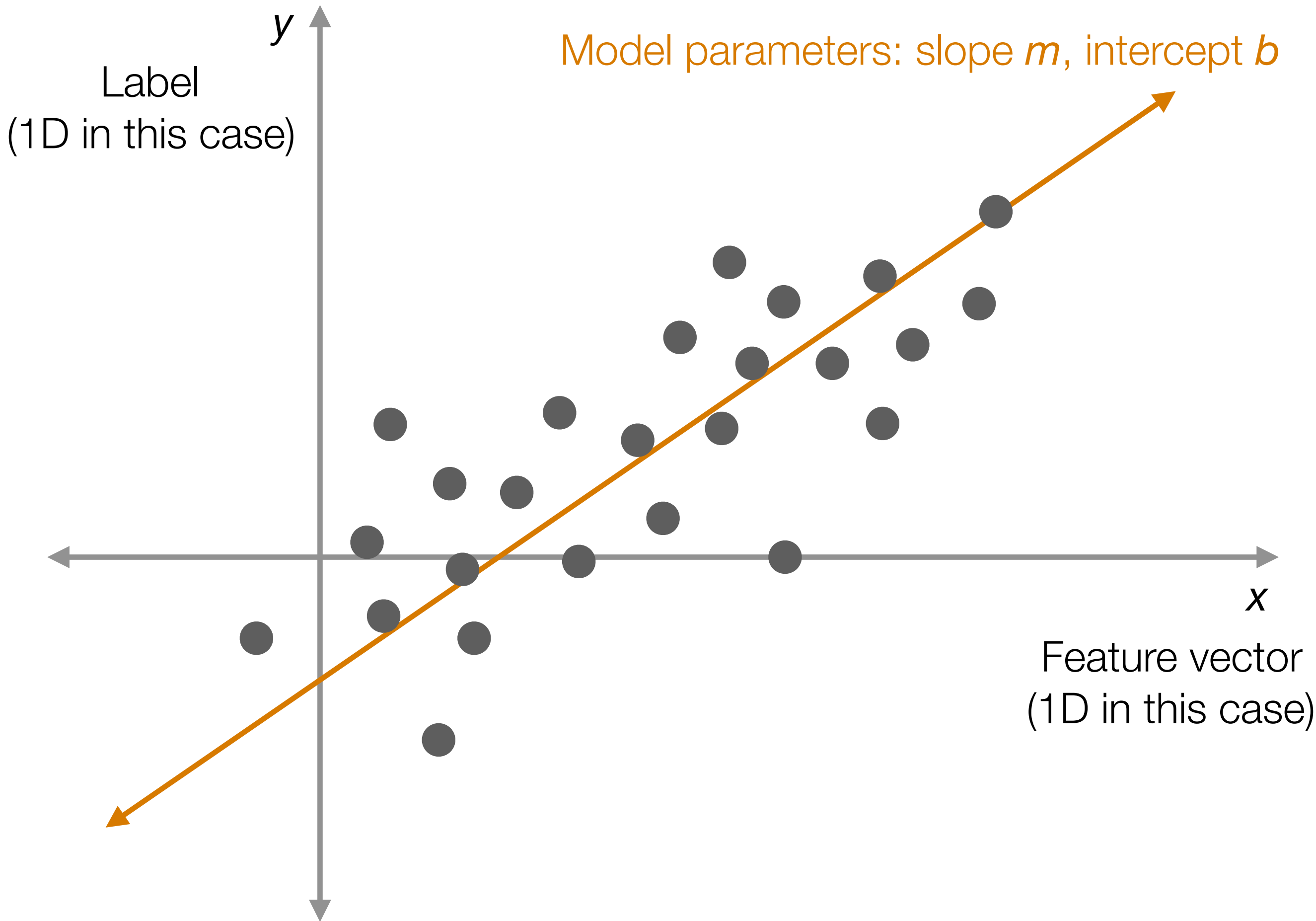
What should the label of this new point be?

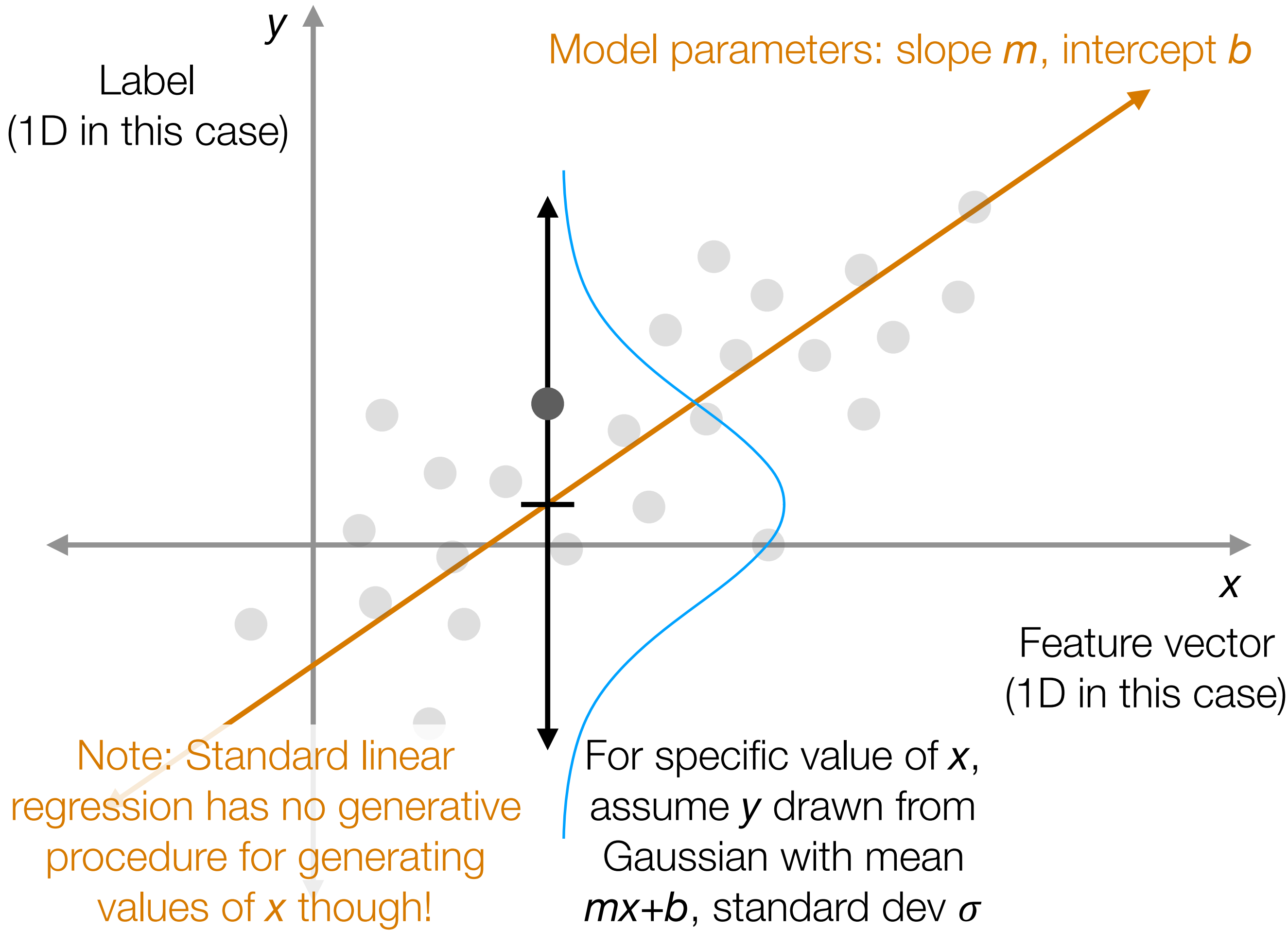
Whichever cluster has higher probability!

This classifier we've created assumes a *generative model*

**You've seen a prediction model
that is partly a generative model**

Linear regression!





Model parameters: slope m , intercept b

Label
(1D in this case)

Feature vector
(1D in this case)

Note: Standard linear regression has no generative procedure for generating values of x though!

For specific value of x , assume y drawn from Gaussian with mean $mx+b$, standard dev σ

Predictive Data Analysis

Training data

$$(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$$

Goal: Given new feature vector x , predict label y

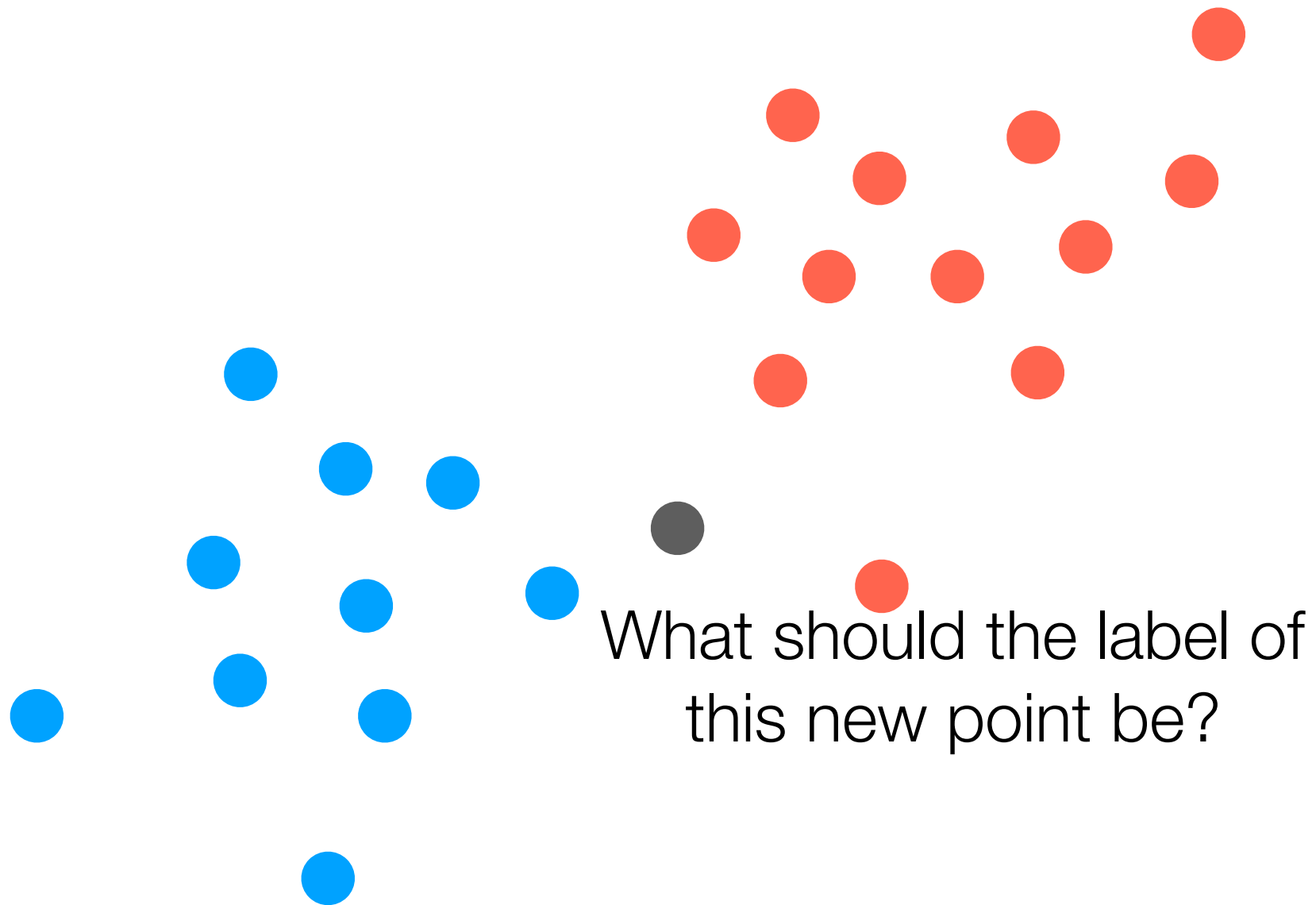
- y is discrete (such as colors **red** and **blue**)
→ prediction method is called a **classifier**
- y is continuous (such as a real number)
→ prediction method is called a **regressor**

A giant zoo of methods

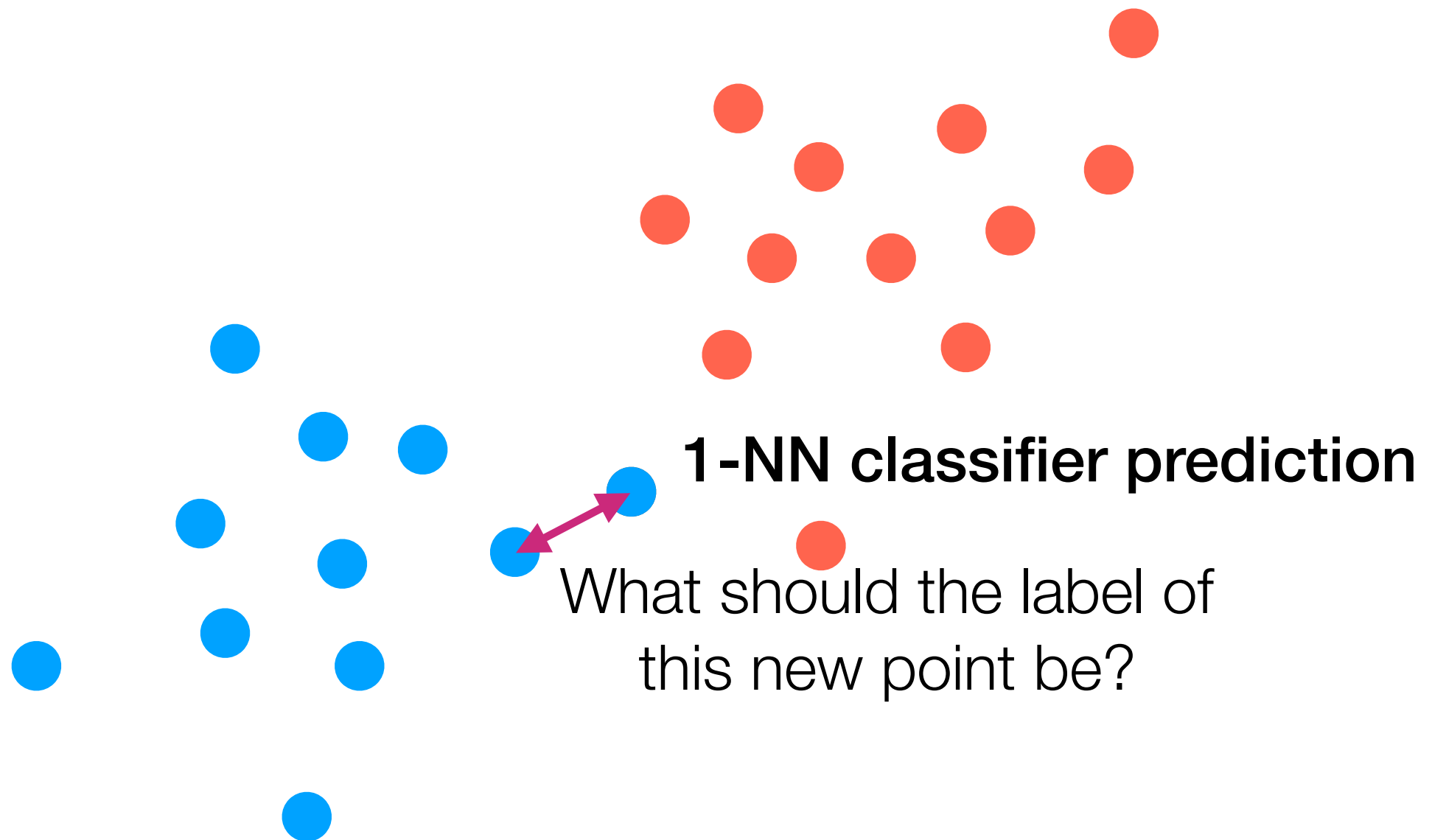
- Generative models (like what we just described)
- Discriminative methods (just care about learning prediction rule *without* assuming generative model)

Example of a Discriminative Method: k -NN Classification

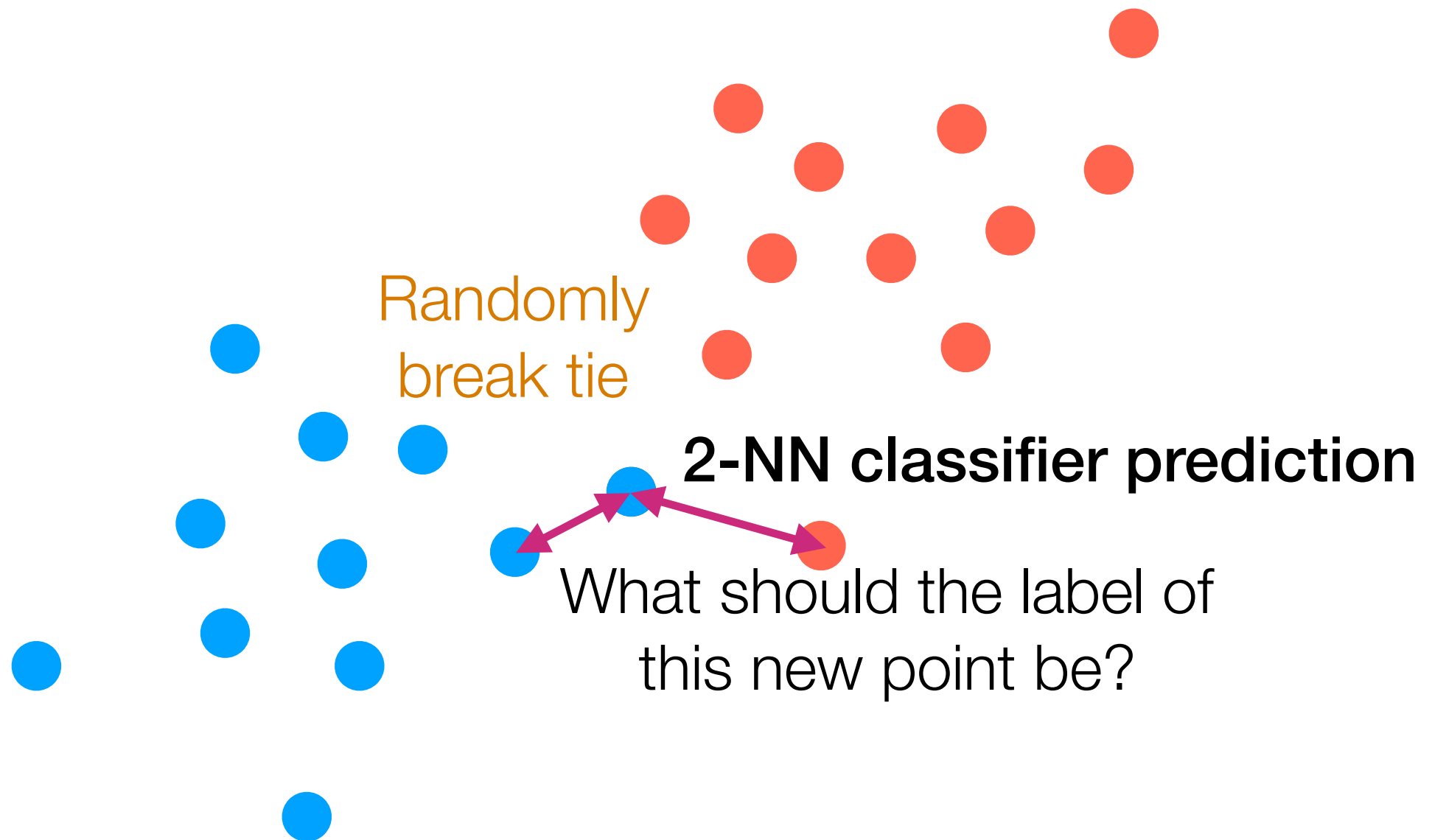
Example: k -NN Classification



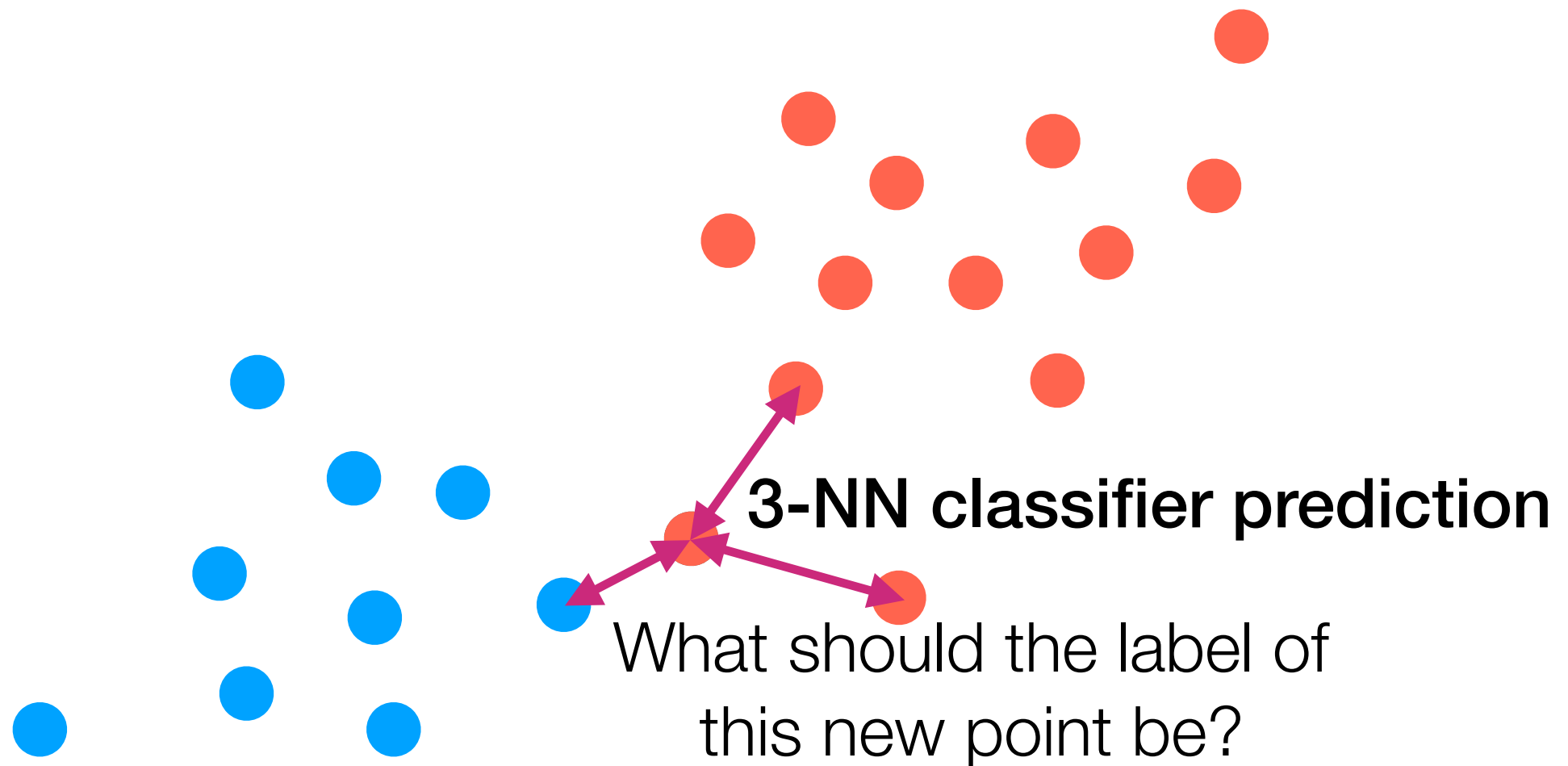
Example: k -NN Classification



Example: k -NN Classification



Example: k -NN Classification



● We just saw: $k = 1$, $k = 2$, $k = 3$

What happens if $k = n$?

How do we choose k ?

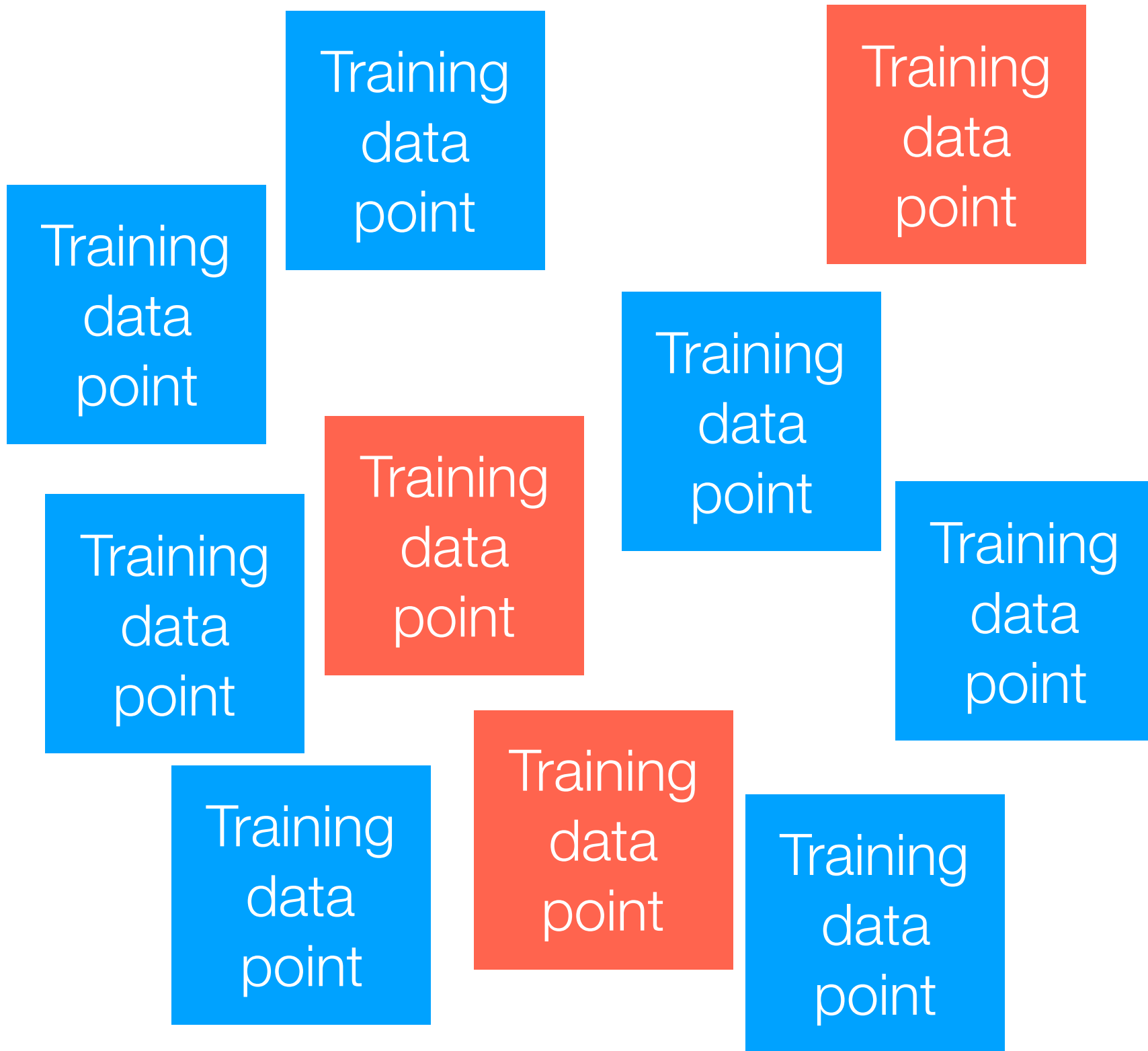
What I'll describe next can be used to select hyperparameter(s) for any prediction method

First: How do we assess how good a prediction method is?

Hyperparameters vs. Parameters

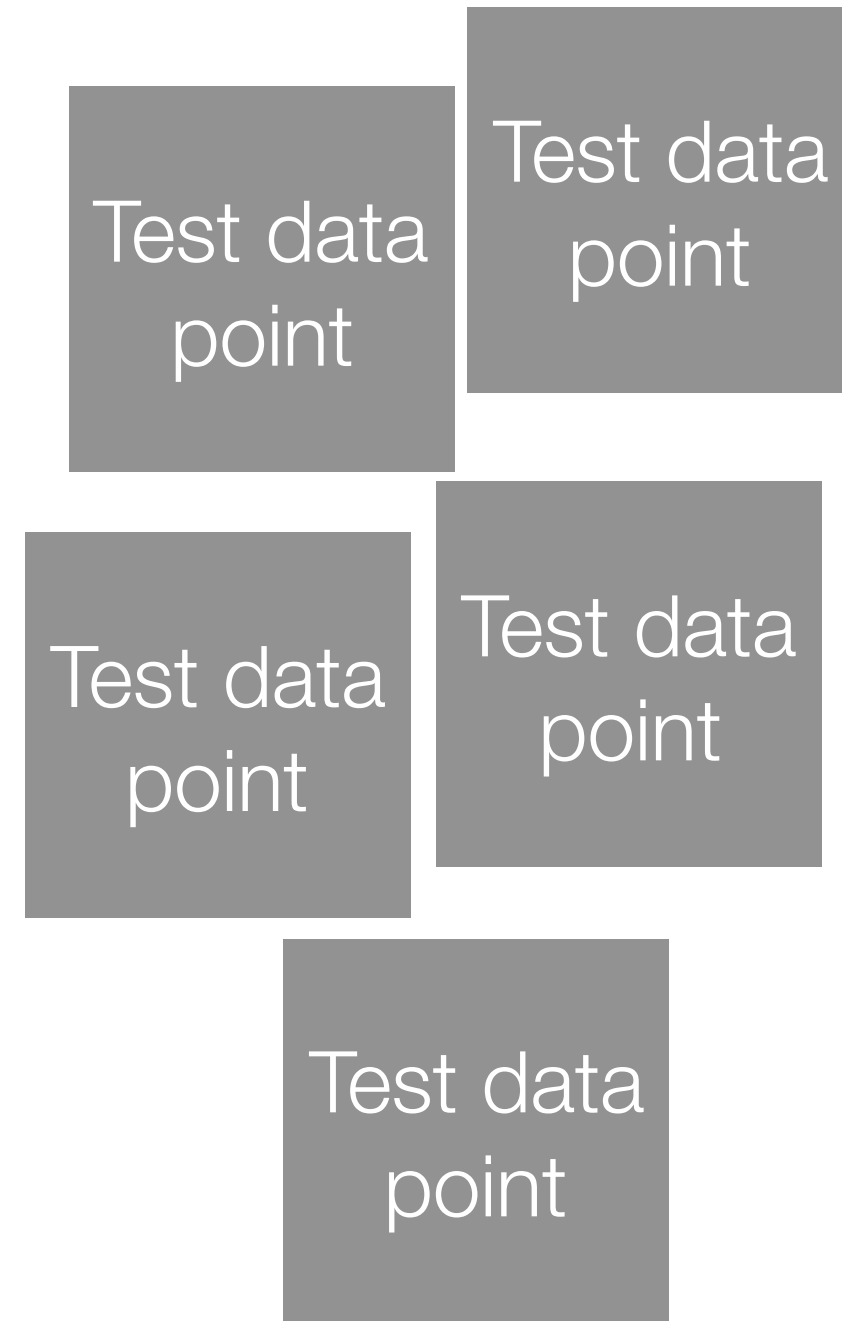
- We fit a model's parameters to training data (terminology: we “learn” the parameters)
- We pick values of hyperparameters and they do *not* get fit to training data
- Example: Gaussian mixture model
 - Hyperparameter: number of clusters k
 - Parameters: cluster probabilities, means, covariances
- Example: k -NN classification
 - Hyperparameter: number of nearest neighbors k
 - Parameters: N/A

Training data



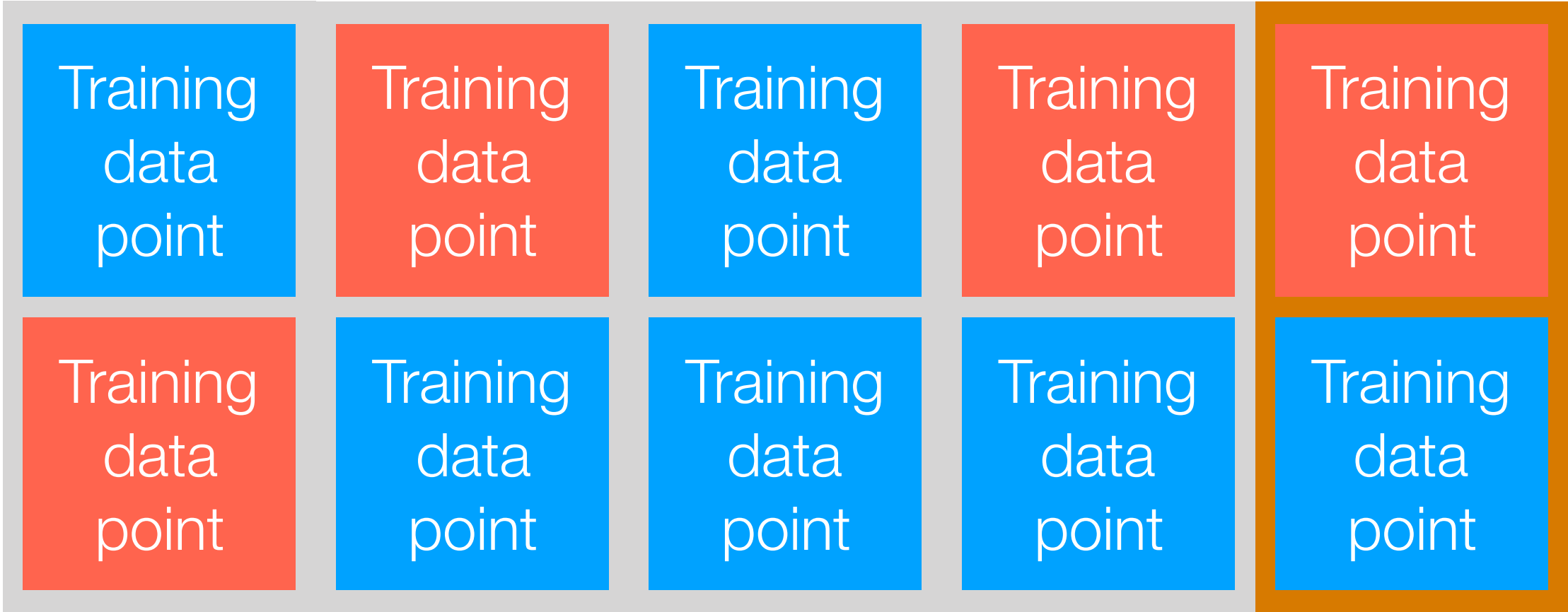
Example: Each data point is an email and we know whether it is spam/ham

Want to classify these points correctly



Example: future emails to classify as spam/ham

Predicted labels



Train method on data in gray

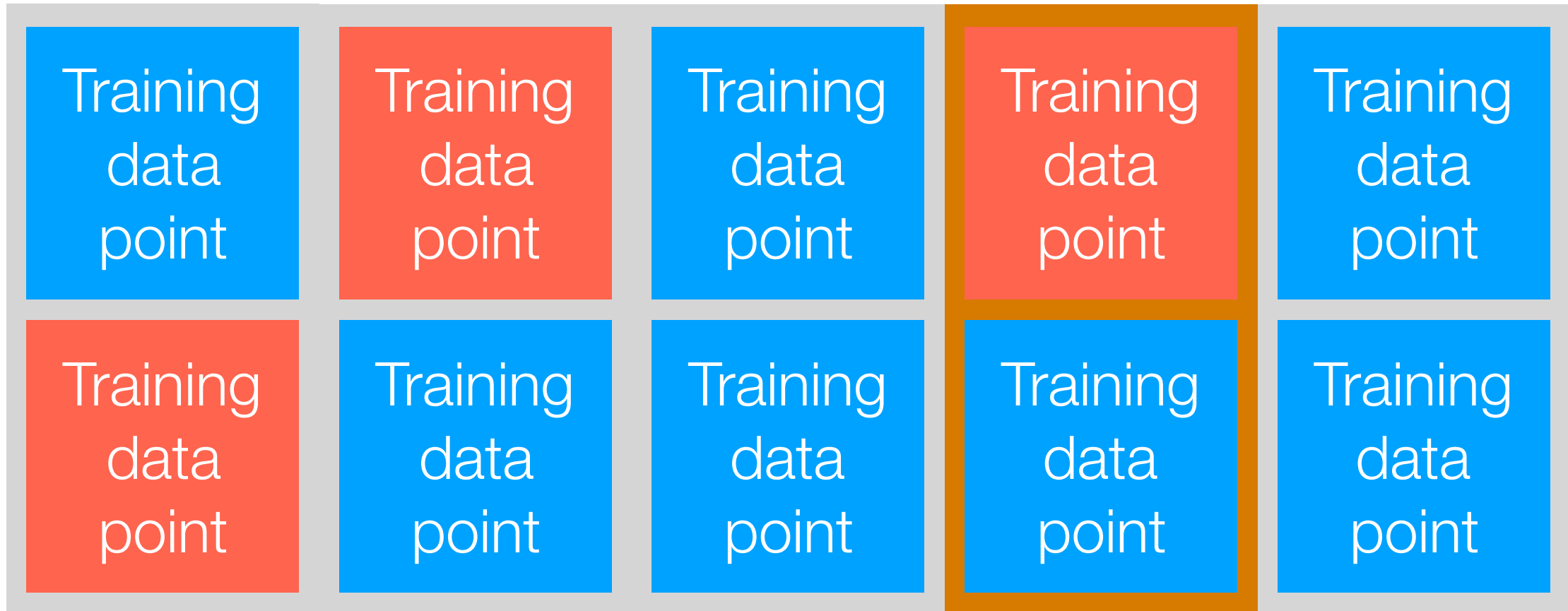
Predict on data in orange

Compute prediction error

Simple data splitting (commonly called train/test split)

50%

In this example: we did a 80%-20% split



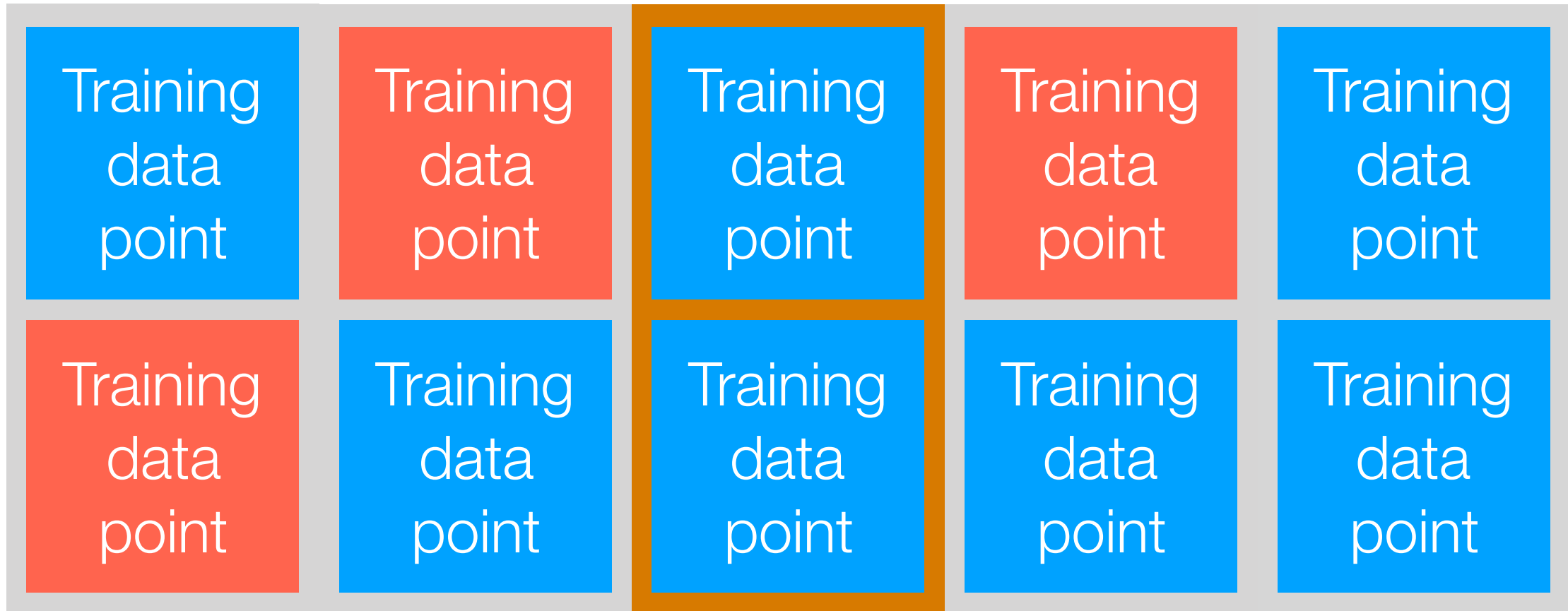
Train method on data in gray

Predict on data in orange

Compute prediction error

0%

50%



Train method on data in gray

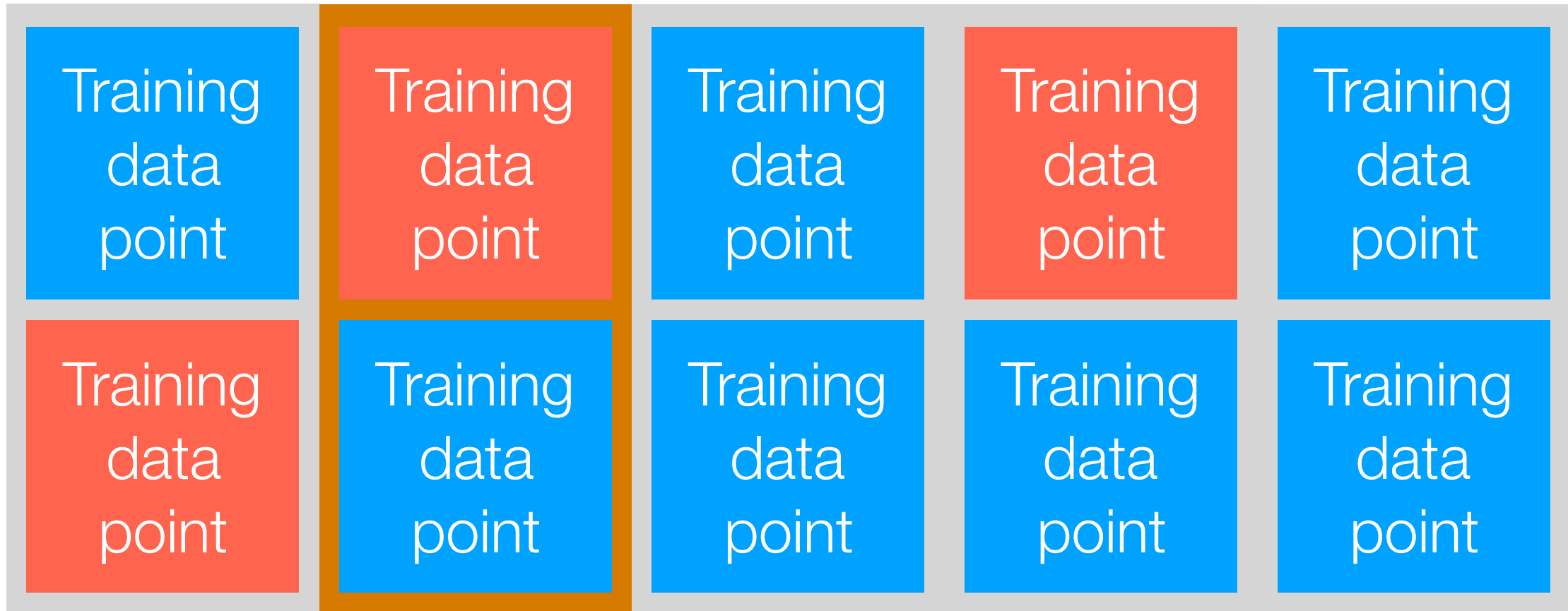
Predict on data
in orange

Compute
prediction error

50%

0%

50%



Train method on data in gray

Predict on data in orange

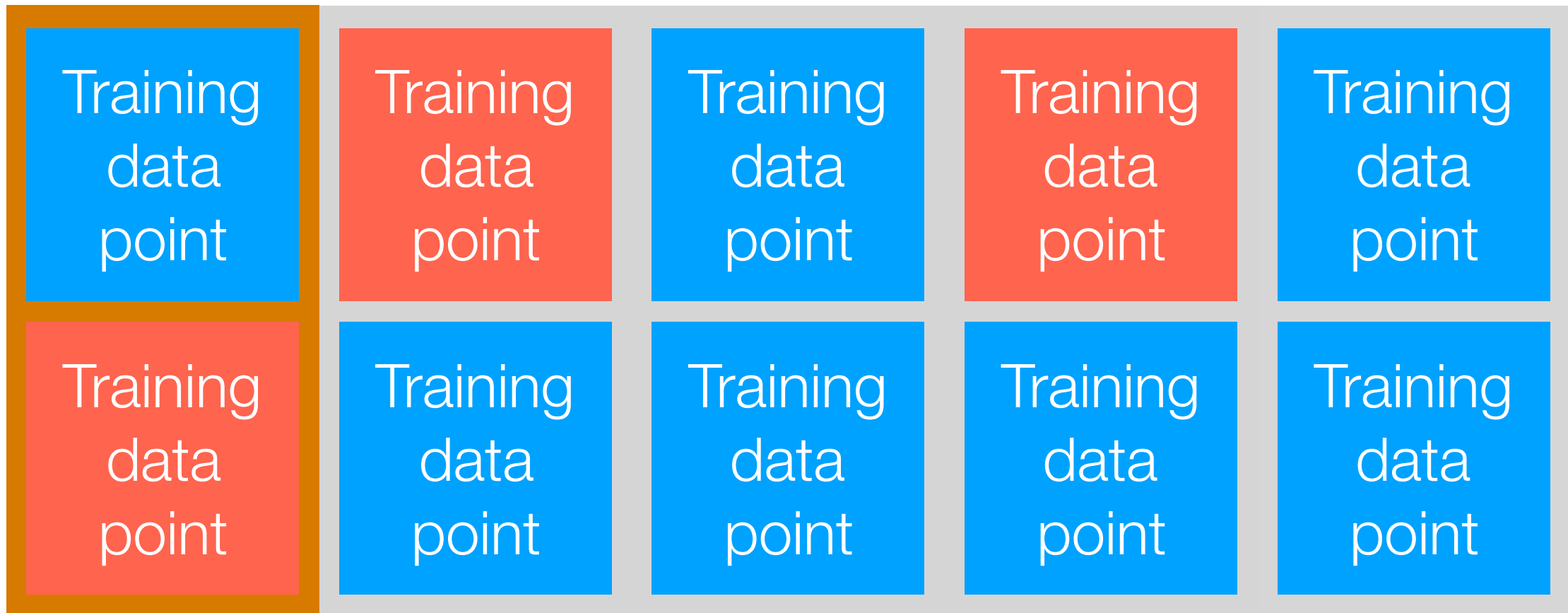
Compute prediction error

0%

50%

0%

50%



Train method on data in gray

Predict on data in orange

Compute prediction error

0%

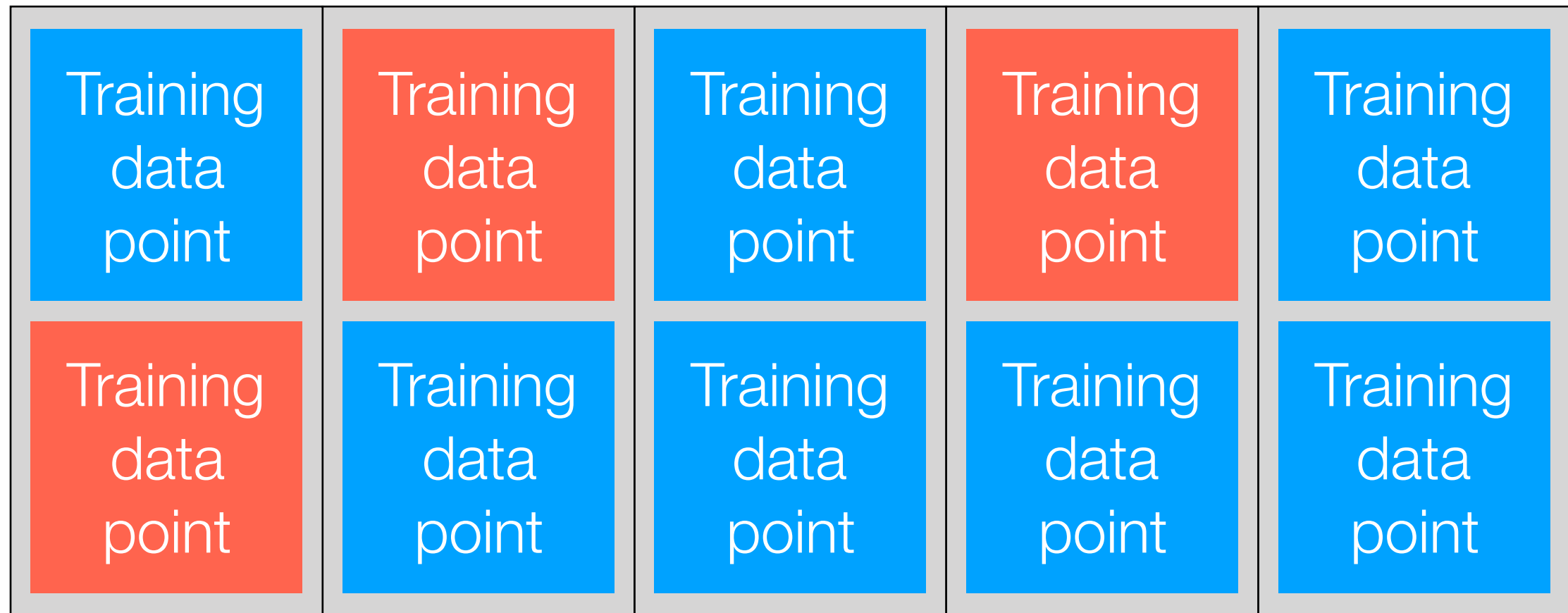
0%

50%

0%

50%

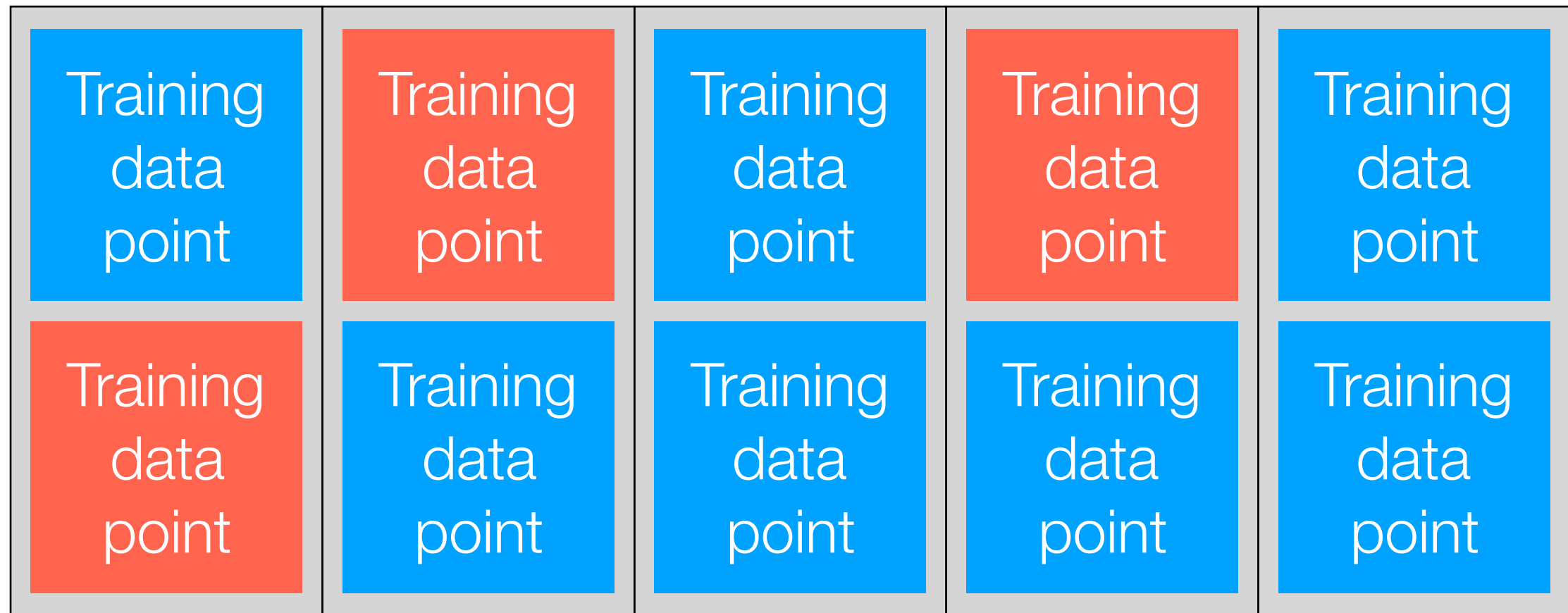
Average error: $(0+0+50+0+50)/5 = 20\%$



1. Shuffle data and put them into “folds” (5 folds in this example)
2. For each fold (which consists of its own train/validation sets):
 - (a) Train on fold’s training data, test on fold’s validation data
 - (b) Compute prediction error
3. Compute average prediction error across the folds

not the same k as in k -means or k -NN classification

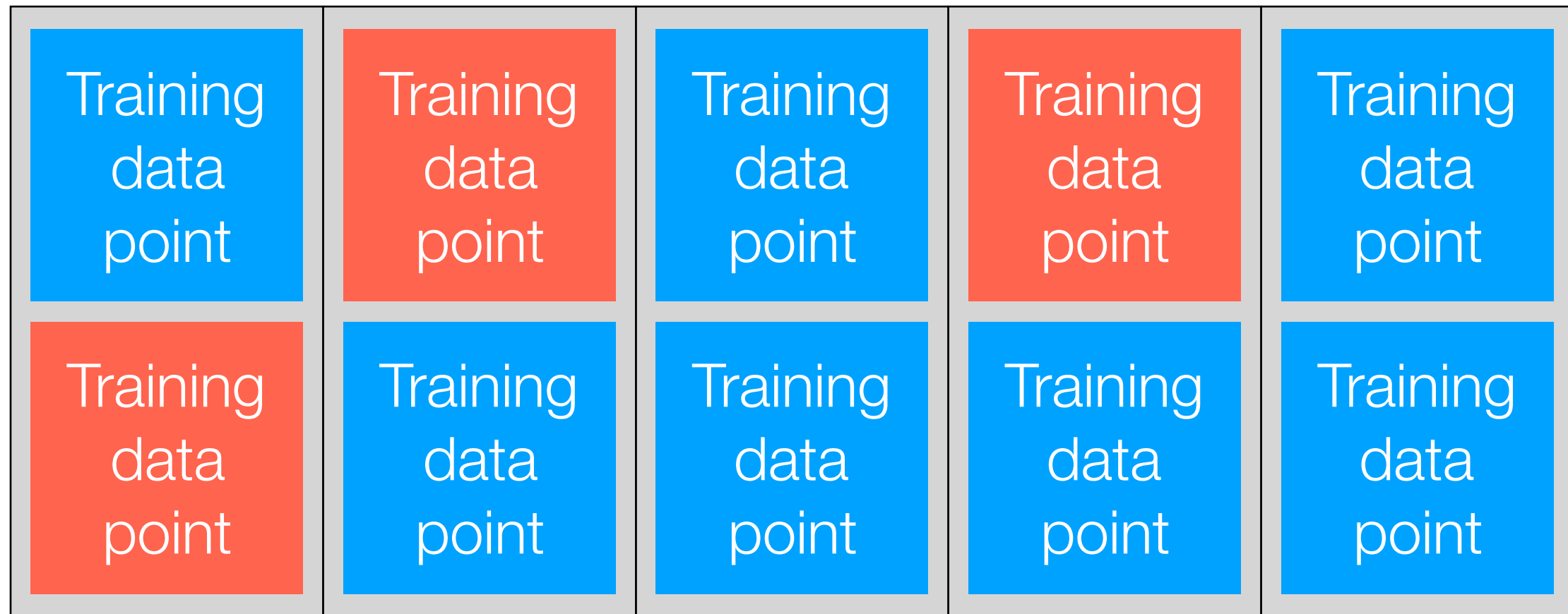
k -fold Cross Validation



1. Shuffle data and put them into “folds” ($k=5$ folds in this example)
2. For each fold (which consists of its own train/validation sets):
 - (a) Train on fold’s training data, test on fold’s validation data
 - (b) Compute prediction error
3. Compute average prediction error across the folds

not the same k as in k -means or k -NN classification

k -fold Cross Validation



1. Shuffle data and put them into “folds” ($k=5$ folds in this example)
2. For each fold (which consists of its own train/validation sets):
 - (a) Train on fold’s training data, test on fold’s validation data
 - (b) Compute **some sort of prediction score**
3. Compute **average prediction score** across the folds
“cross validation score”

Choosing k in k -NN Classification

Note: k -NN classifier has a single hyperparameter k

For each $k = 1, 2, 3, \dots$, the maximum k you are willing to try:

 Compute 5-fold cross validation score using k -NN classifier as prediction method

Use whichever k has the best cross validation score

Automatic Hyperparameter Selection

Suppose the prediction algorithm you're using has hyperparameters θ

For each hyperparameter setting θ you are willing to try:

Compute 5-fold cross validation score using your algorithm with hyperparameters θ

Use whichever θ has the best cross validation score

Why 5?

People have found using 10 folds or 5 folds to work well in practice but it's just empirical — there's no deep reason

Training data

Training data point

Training data point

Important: the errors from simple data splitting and cross-validation are *estimates* of the true error on test data!

Example: earlier, we got a cross validation score of 20% error

This is a guess for the error we will get on test data

This guess is not always accurate!

Example: Each data point is an email and we know whether it is spam/ham

Want to classify these points correctly

Test data point

Test data point

Test data point

Test data point

Test data point

Example: future emails to classify as spam/ham

Cross-Validation Remarks

- k -fold cross-validation is a randomized procedure
 - Re-running CV results in different cross-validation scores!
- Suppose there are n training data points and k folds
 - If we are trying 10 different hyperparameter settings, how many times do we do model fitting? $10k$
 - If this number is similar in size to n , CV can overfit!
 - How many training data are used in each model fit during cross-validation? $[(k-1)/k]n$
 - Smaller # folds typically means faster training
- If $k = n$, would re-running cross-validation result in different cross-validation scores? What about $k = 2$?
 - For deterministic training procedure: same CV result for $k = n$ (since shuffling doesn't matter), different for $k = 2$

Different Ways to Measure Accuracy

Simplest way:

- **Raw error rate:** fraction of predicted labels that are wrong (this was in our cross validation example earlier)

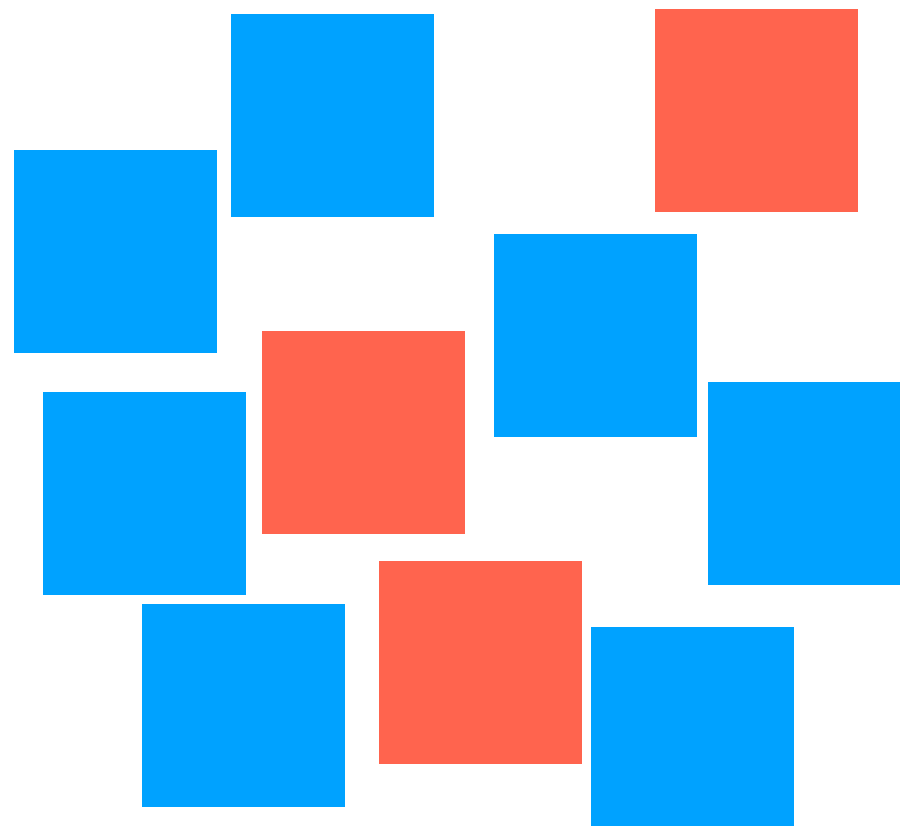
In “binary” classification (there are 2 labels such as spam/ham) when 1 label is considered “positive” and the other “negative”:

Different Ways to Measure Accuracy

Simplest way:

- **Raw error rate:** fraction of predicted labels that are wrong (this was in our cross validation example earlier)

In “binary” classification (there are 2 labels such as spam/ham) when 1 label is considered “positive” and the other “negative”:

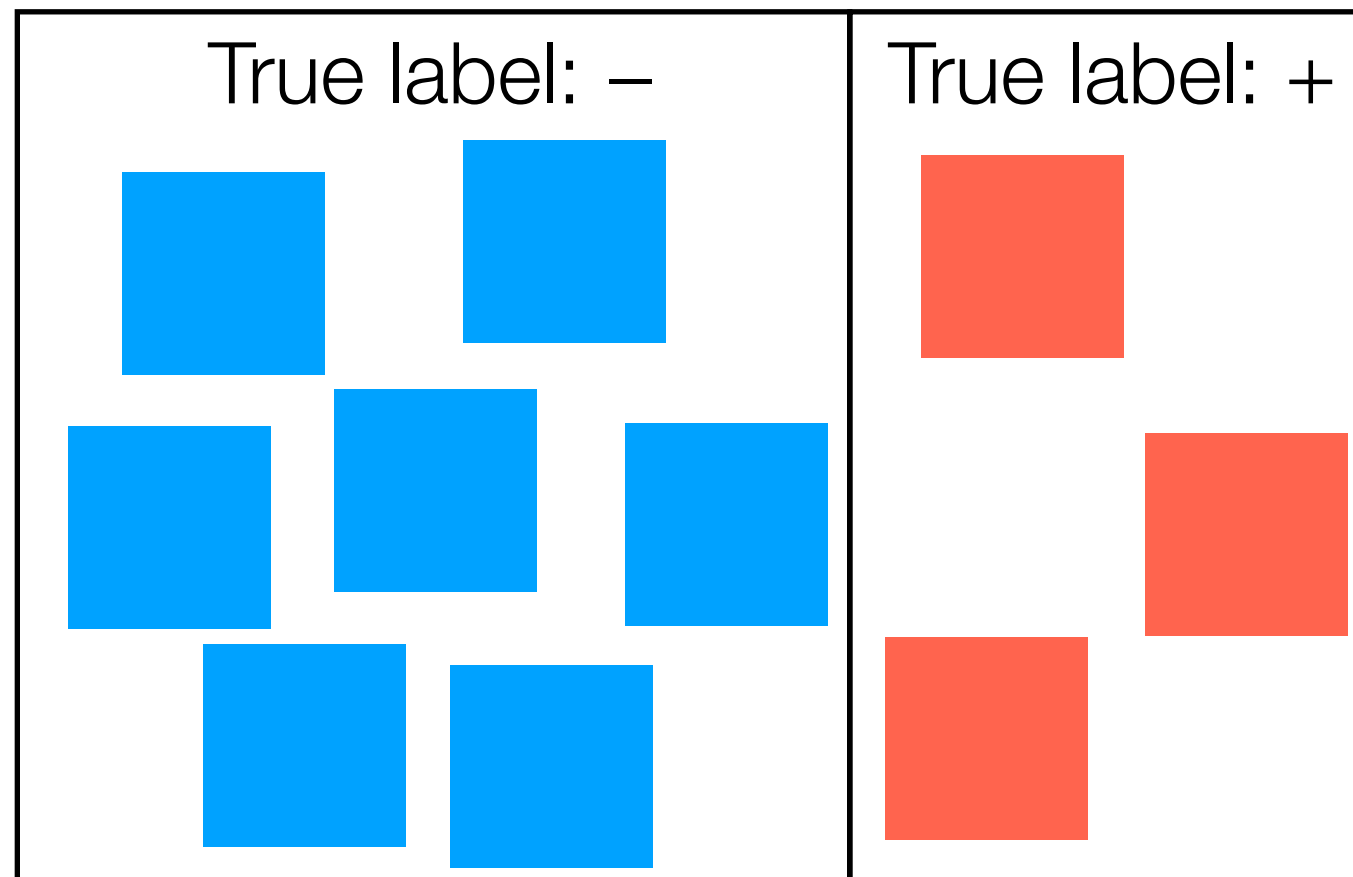


Different Ways to Measure Accuracy

Simplest way:

- **Raw error rate:** fraction of predicted labels that are wrong (this was in our cross validation example earlier)

In “binary” classification (there are 2 labels such as spam/ham) when 1 label is considered “**positive**” and the other “**negative**”:



Different Ways to Measure Accuracy

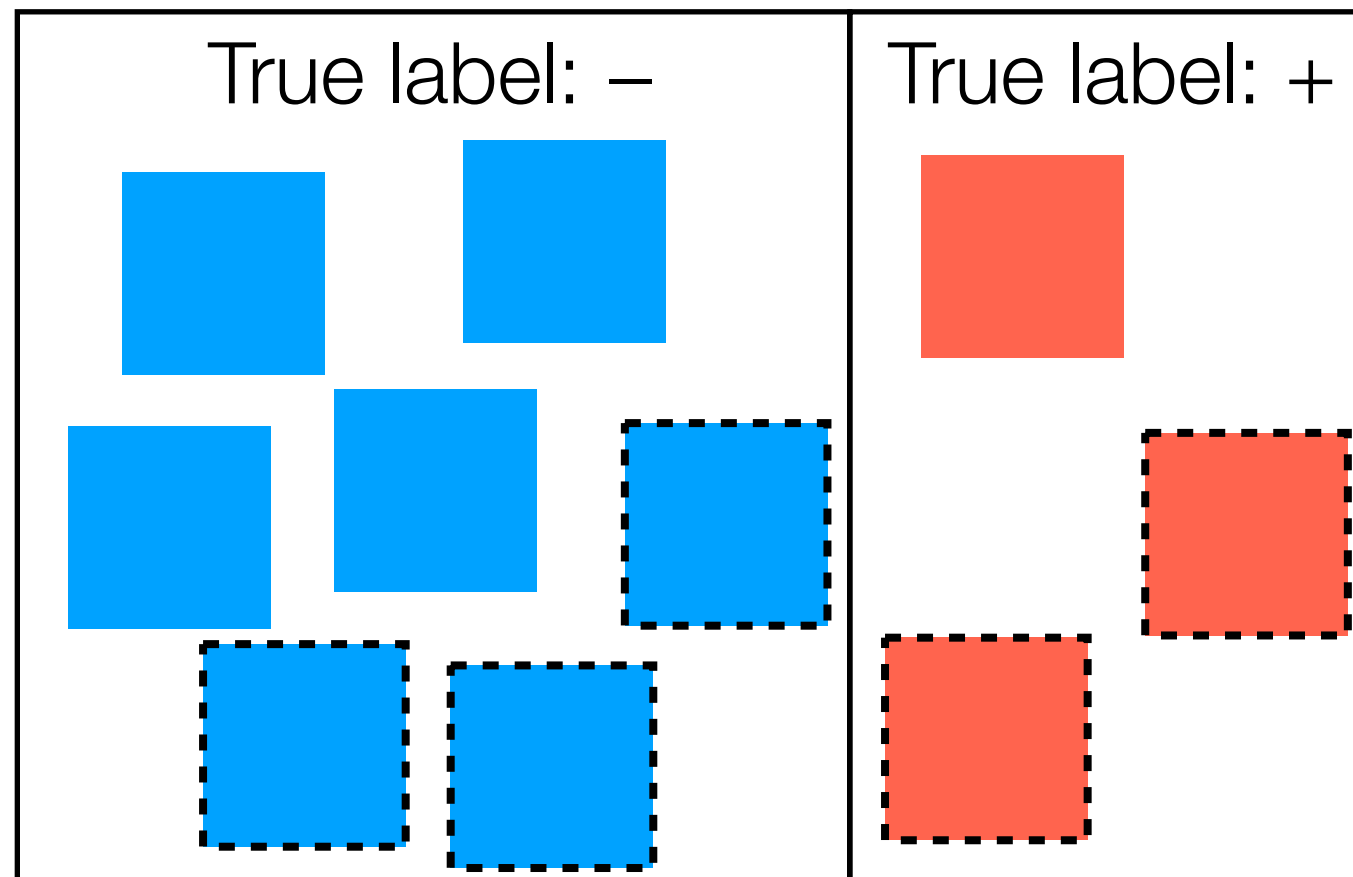
Simplest way:

- **Raw error rate:** fraction of predicted labels that are wrong (this was in our cross validation example earlier)

In “binary” classification (there are 2 labels such as spam/ham) when 1 label is considered “positive” and the other “negative”:

Outlined in dotted black:
predicted label +

(all other points predicted to be -)

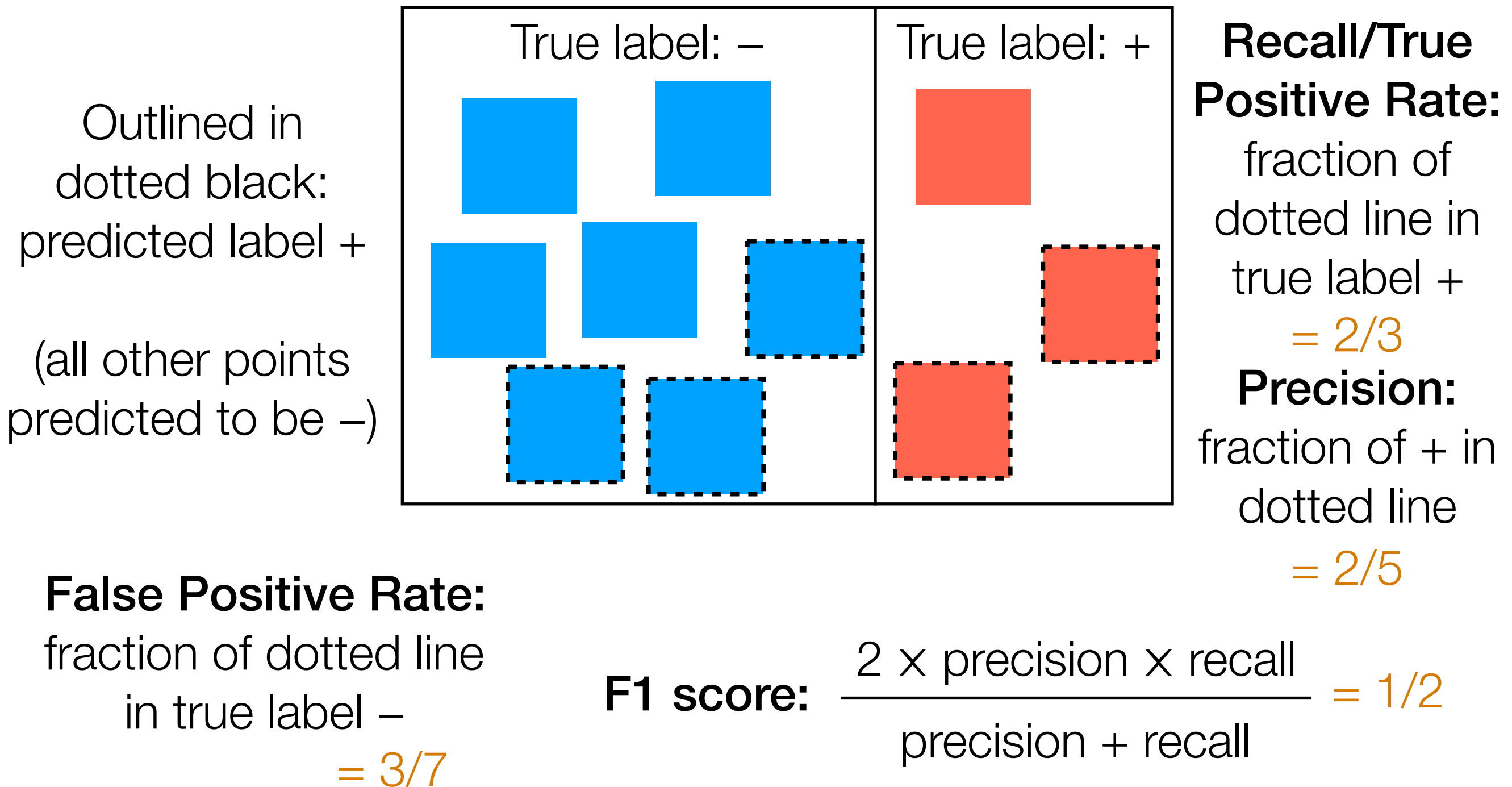


Recall/True Positive Rate:
fraction of dotted line in true label +

Precision:
fraction of + in dotted line

- **Raw error rate:** fraction of predicted labels that are wrong (this was in our cross validation example earlier)

In “binary” classification (there are 2 labels such as spam/ham) when 1 label is considered “positive” and the other “negative”:



Prediction and Model Validation

Demo

Deep Learning

IMAGENET

Over 10 million images, 1000 object classes



2011: Traditional computer vision achieves accuracy ~74%

2012: Initial deep neural network approach accuracy ~84%

2015 onwards: Deep learning achieves accuracy 96%+

Russakovsky et al. ImageNet Large Scale Visual Recognition Challenge. IJCV 2015.

Deep Learning Takeover

Academia:

- Top computer vision conferences (CVPR, ICCV, ECCV) are now nearly all about deep learning
- Top machine learning conferences (ICML, NeurIPS) have *heavily* been taken over by deep learning

Heavily dominated by industry now!

Extremely useful in practice:

- Near human level image classification (including handwritten digit recognition)
- Near human level speech recognition
- Improvements in machine translation, text-to-speech
- Self-driving cars
- *Better* than humans at playing Go





Google DeepMind's AlphaGo vs Lee Sedol, 2016

GAMING

TECH

ARTIFICIAL INTELLIGENCE

DeepMind's StarCraft 2 AI is now better than 99.8 percent of all human players

16 💬

AlphaStar is now grandmaster level in the real-time strategy game

By [Nick Statt](#) | [@nickstatt](#) | Oct 30, 2019, 2:00pm EDT



SHARE



Is it all hype?

Should you as a human be afraid of robots taking your job?!?



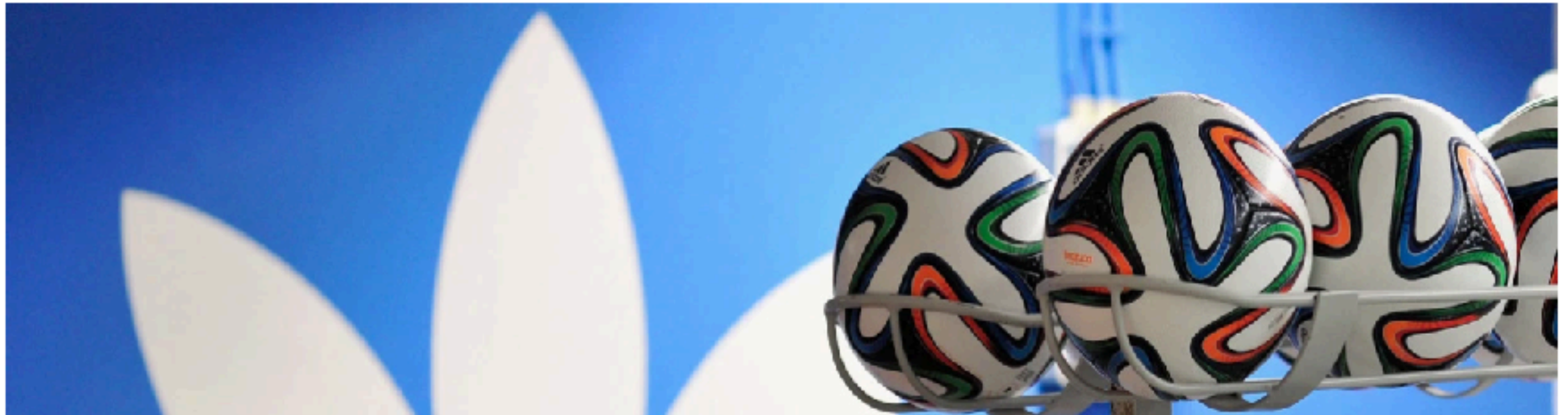
BUSINESS



American robots lose jobs to Asian robots as Adidas shifts manufacturing

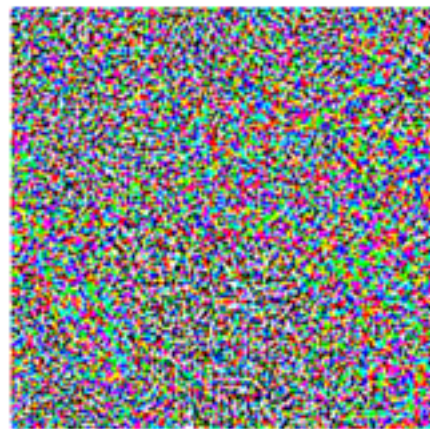
By Reuters

November 11, 2019 | 9:13am | Updated





+ .007 ×



=

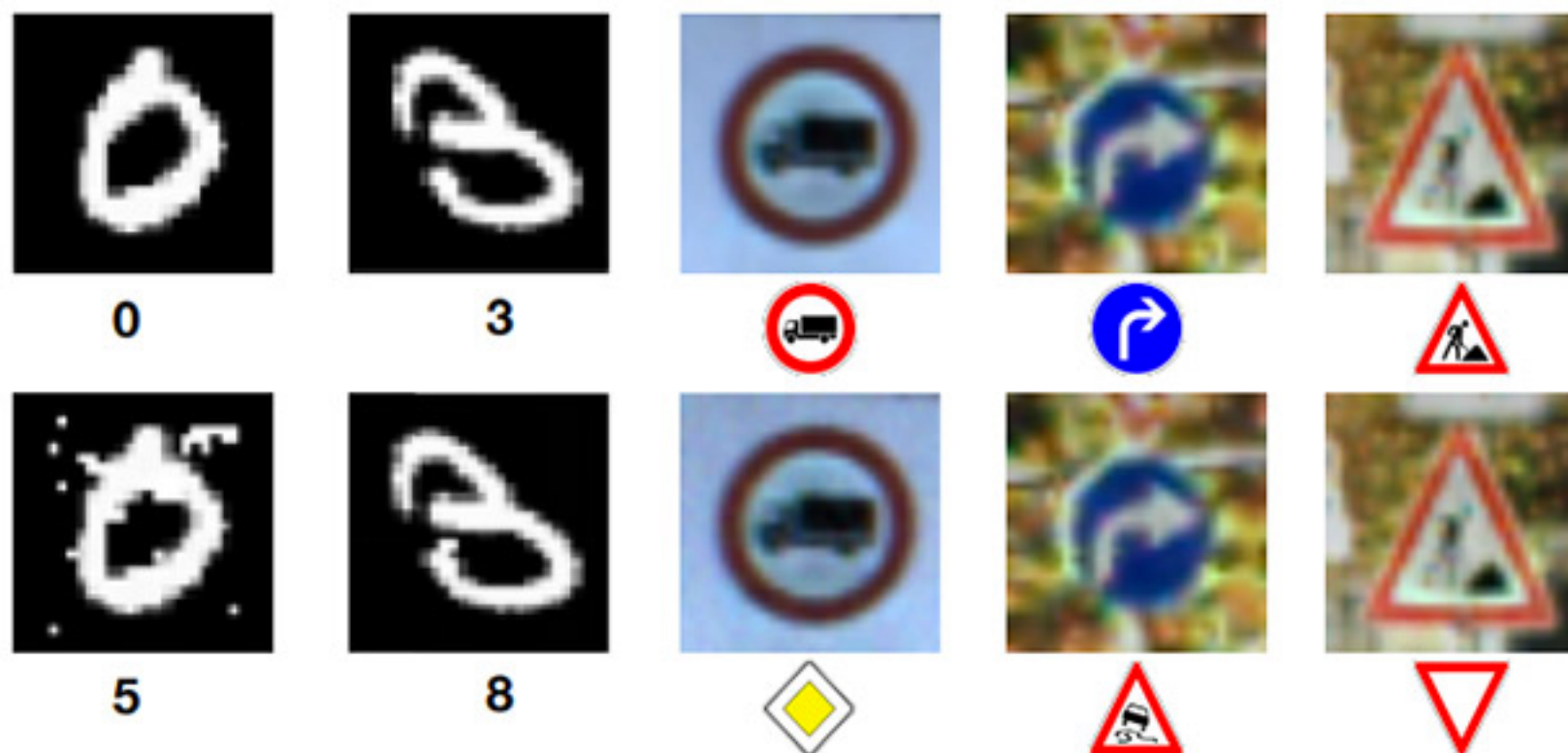


panda
~58% confidence

adversarial
noise

gibbon
~99% confidence

Source: Goodfellow, Shlens, and Szegedy. Explaining and Harnessing Adversarial Examples. ICLR 2015.

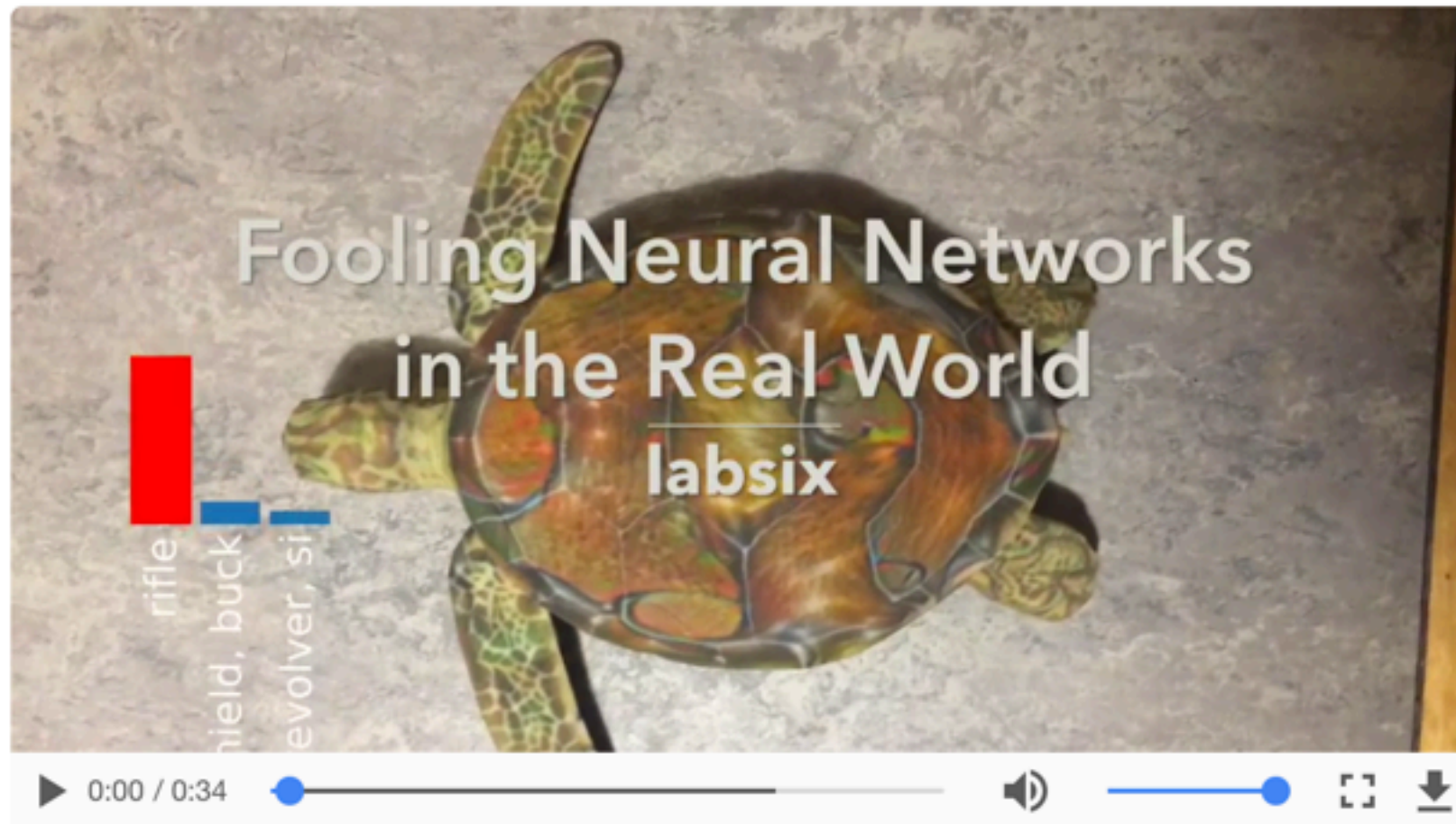


Source: Papernot et al. Practical Black-Box Attacks against Machine Learning. Asia Conference on Computer and Communications Security 2017.

Fooling Neural Networks in the Physical World with 3D Adversarial Objects

31 Oct 2017 · 3 min read — shared on [Hacker News](#), [Lobsters](#), [Reddit](#), [Twitter](#)

We've developed an approach to generate *3D adversarial objects* that reliably fool neural networks in the real world, no matter how the objects are looked at.



Neural network based classifiers reach near-human performance in many tasks, and they're used in high risk, real world systems. Yet, these same neural networks are particularly vulnerable to *adversarial examples*, carefully perturbed inputs that cause

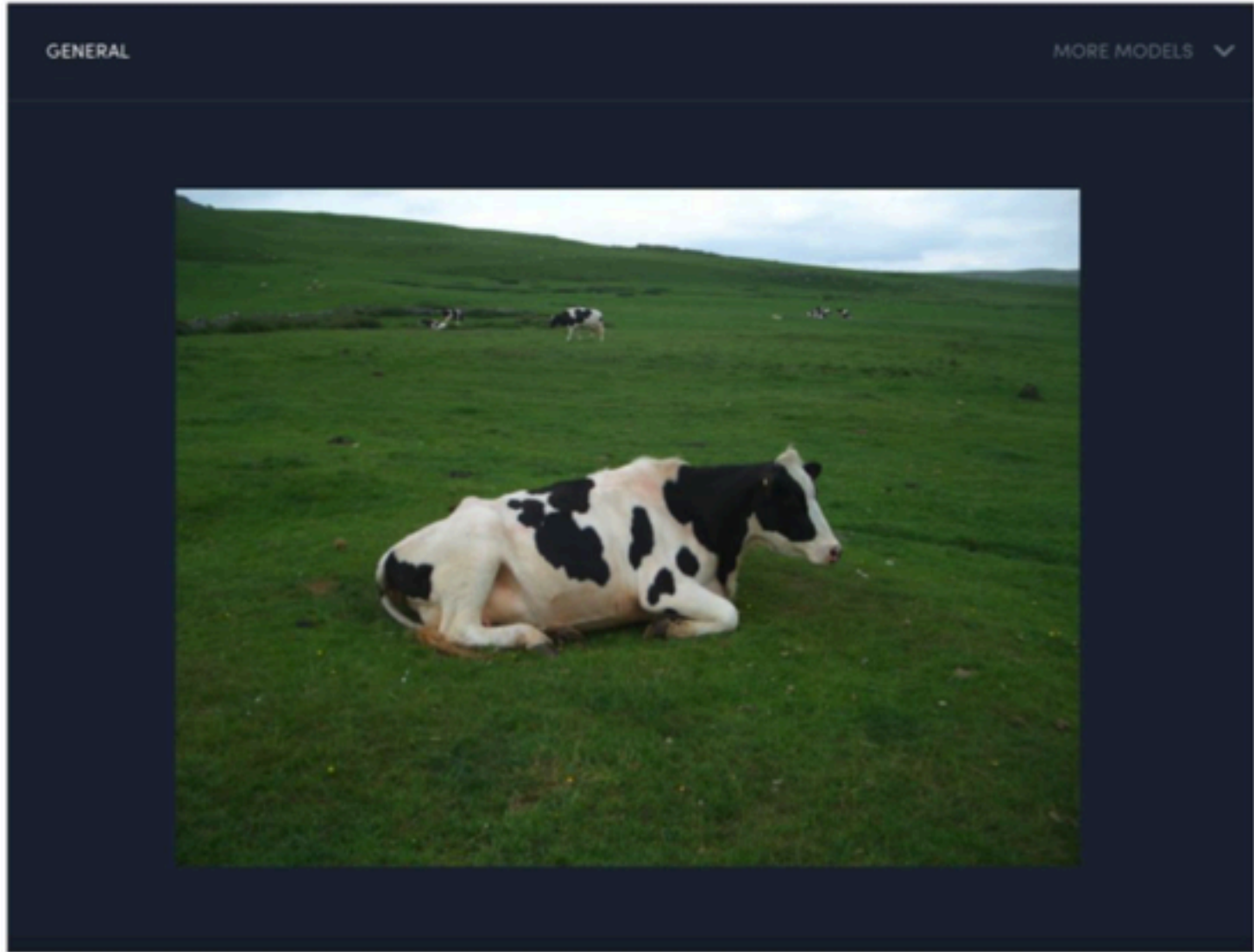
Source: labsix



Source: <https://www.cc.gatech.edu/news/611783/erasing-stop-signs-shapeshifter-shows-self-driving-cars-can-still-be-manipulated>

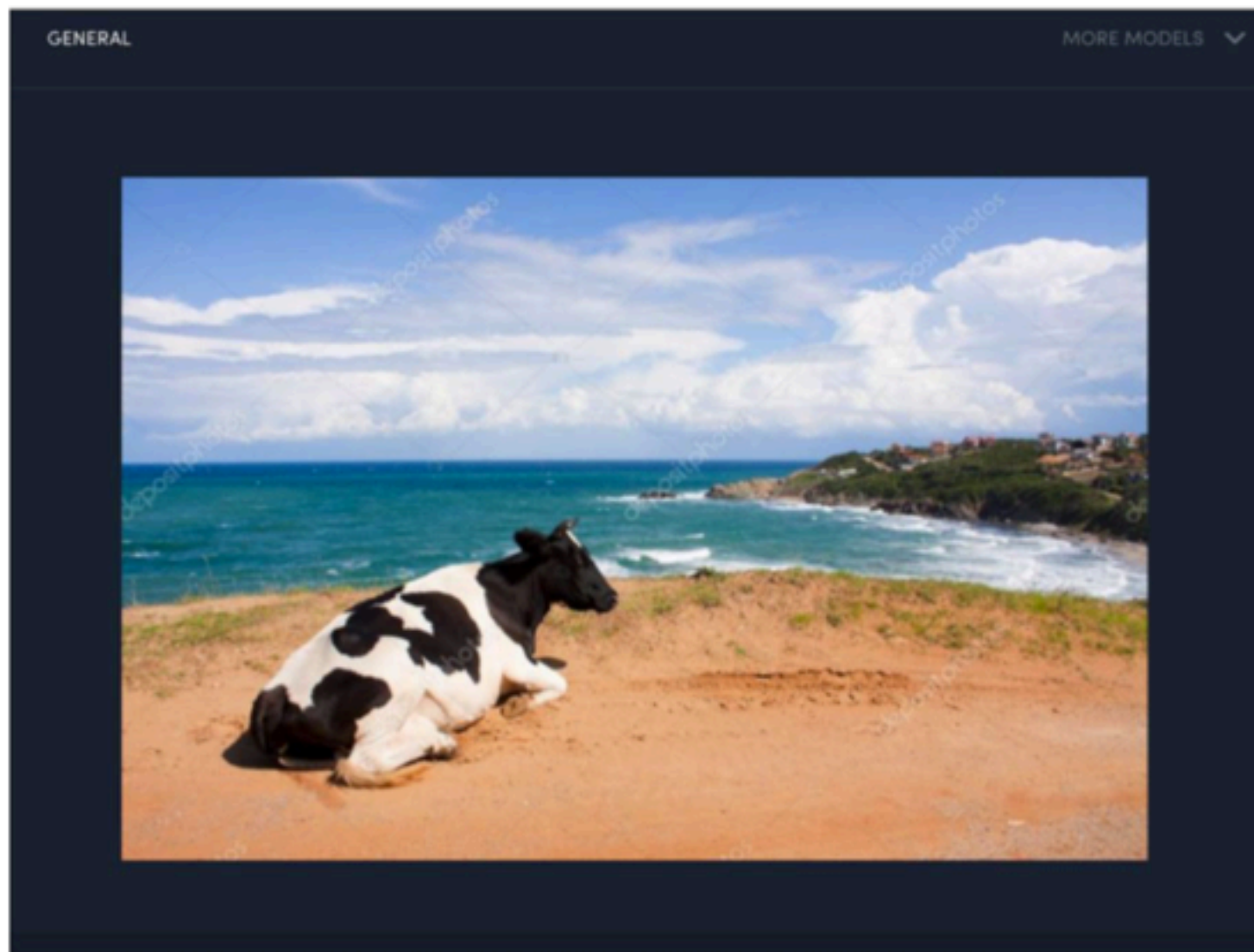


Source: Gizmodo article "This Neural Network's Hilariously Bad Image Descriptions Are Still Advanced AI". September 16, 2015. (They're using the NeuralTalk image-to-caption software.)



General	VIEW DOCS
cow	0.992
cattle	0.983
mammal	0.979
grass	0.978
livestock	0.966
farm	0.964
landscape	0.963
pasture	0.954
grassland	0.949
agriculture	0.948
no person	0.945

Source: Pietro Perona




General [VIEW DOCS](#)

no person	0.991
beach	0.990
water	0.985
sand	0.981
sea	0.980
travel	0.978
seashore	0.972
summer	0.954
sky	0.946
outdoors	0.944
ocean	0.936

cow is not among top objects found!

Source: Pietro Perona

GENERAL FACE NSFW COLOR MORE MODELS



PREDICTED CONCEPT	PROBABILITY
group	0.979
adult	0.977
people	0.976
furniture	0.960
room	0.957
business	0.903
indoors	0.901
man	0.896
seat	0.895

VIEW DOCS

elephant is not among top objects found!

Source: David Lopez-Paz

Another AI Winter?

~1970's: First AI winter over symbolic AI

~1980's: Second AI winter over "expert systems"

Every time: Lots of hype, explosion in funding, then bubble bursts



Michael Jordan [Follow](#)

Michael I. Jordan is a Professor in the Department of Electrical Engineering and Computer Sciences and the Department of Statistics at UC Berkeley.

Apr 18 · 16 min read

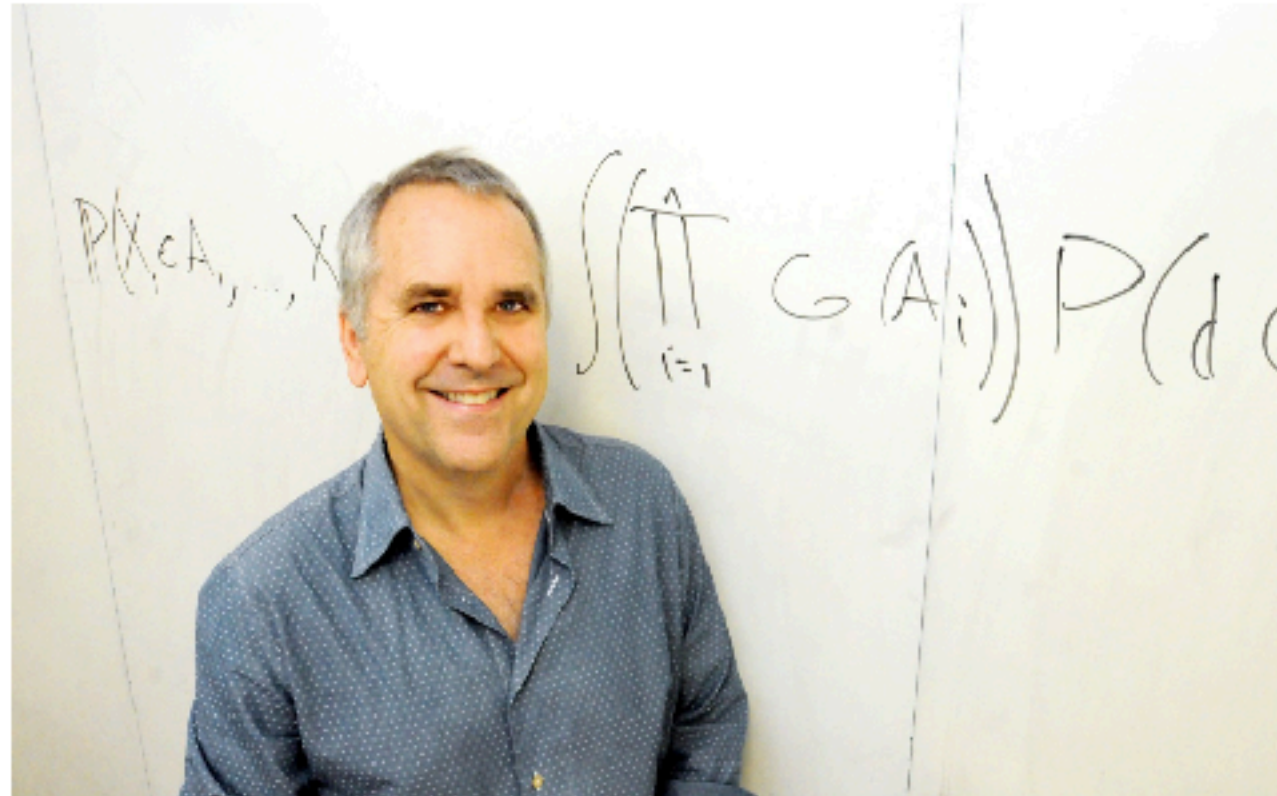


Photo credit: Peg Skorpinski

Artificial Intelligence—The Revolution Hasn't Happened Yet

Artificial Intelligence (AI) is the mantra of the current era. The phrase is intoned by technologists, academicians, journalists and venture capitalists

<https://medium.com/@mijordan3/artificial-intelligence-the-revolution-hasnt-happened-yet-5e1d5812e1e7>

What is deep learning?



Classification units



PIT/AIT



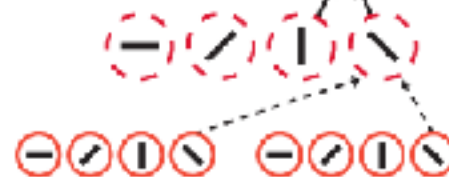
V4/PIT



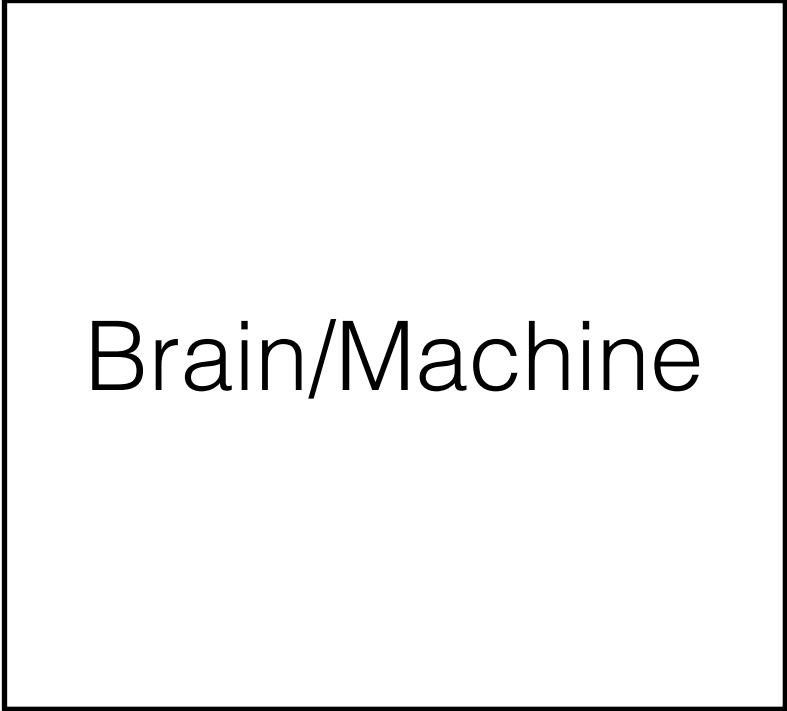
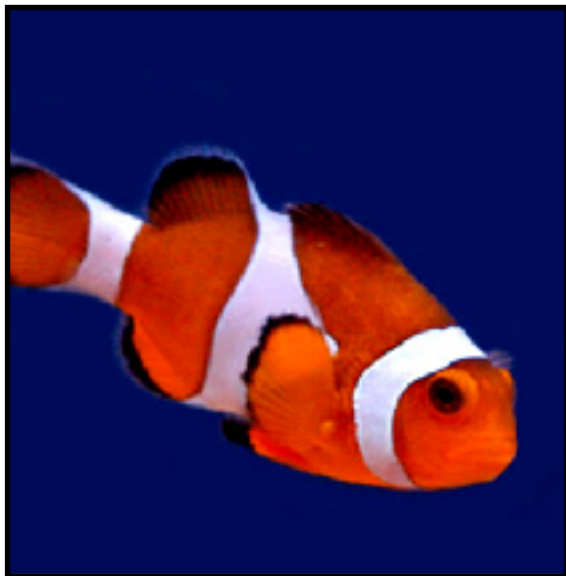
V2/V4



V1/V2

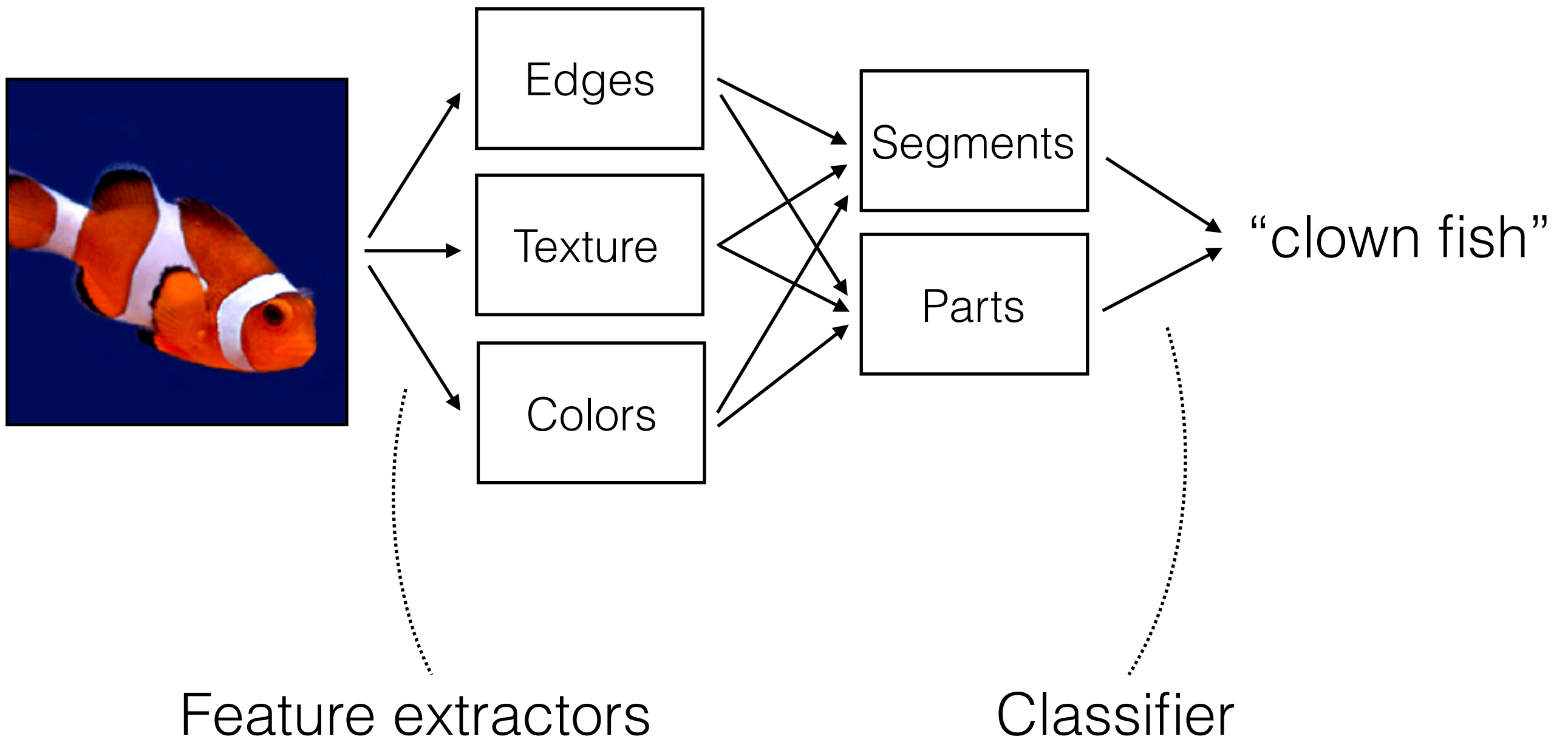


Basic Idea



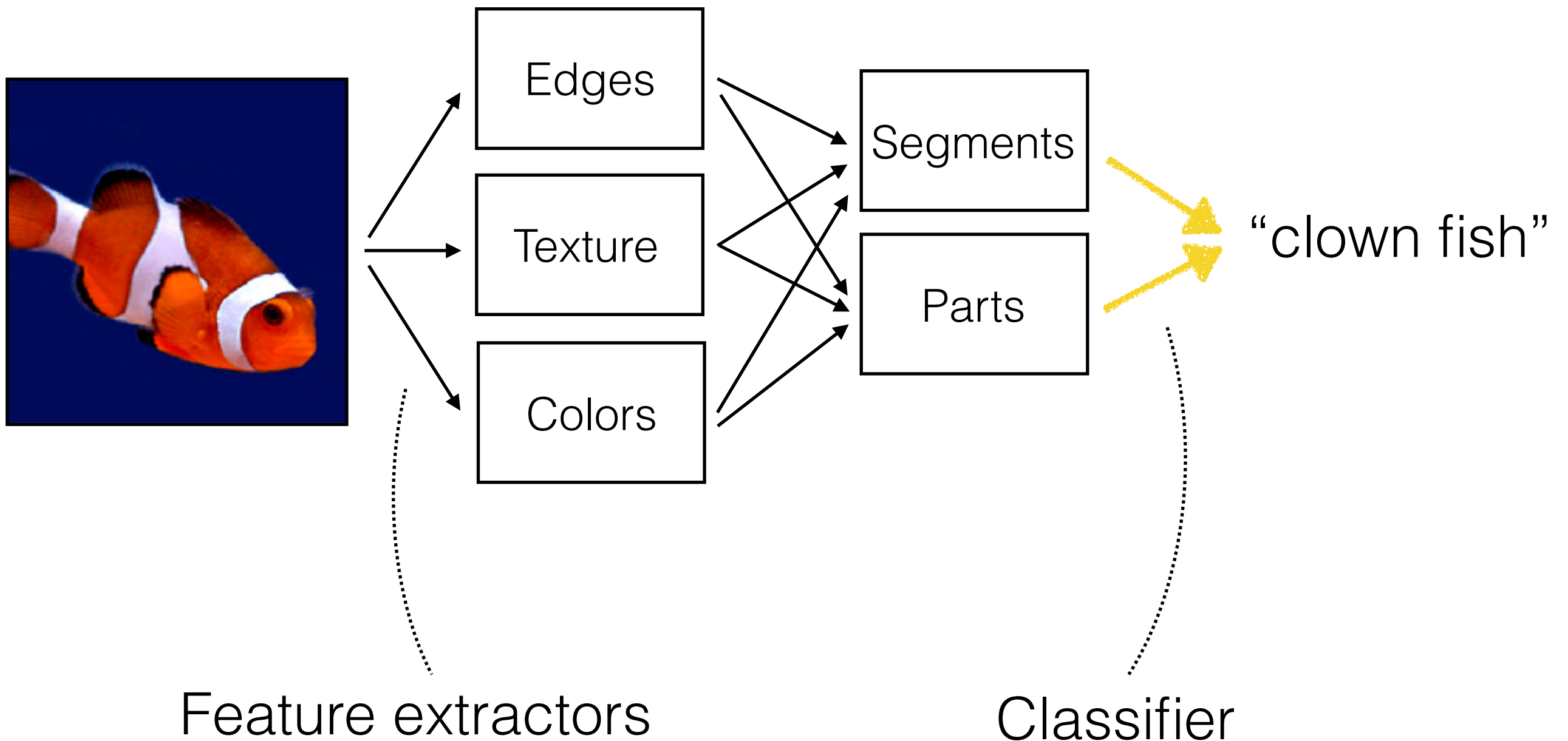
“clown fish”

Object Recognition



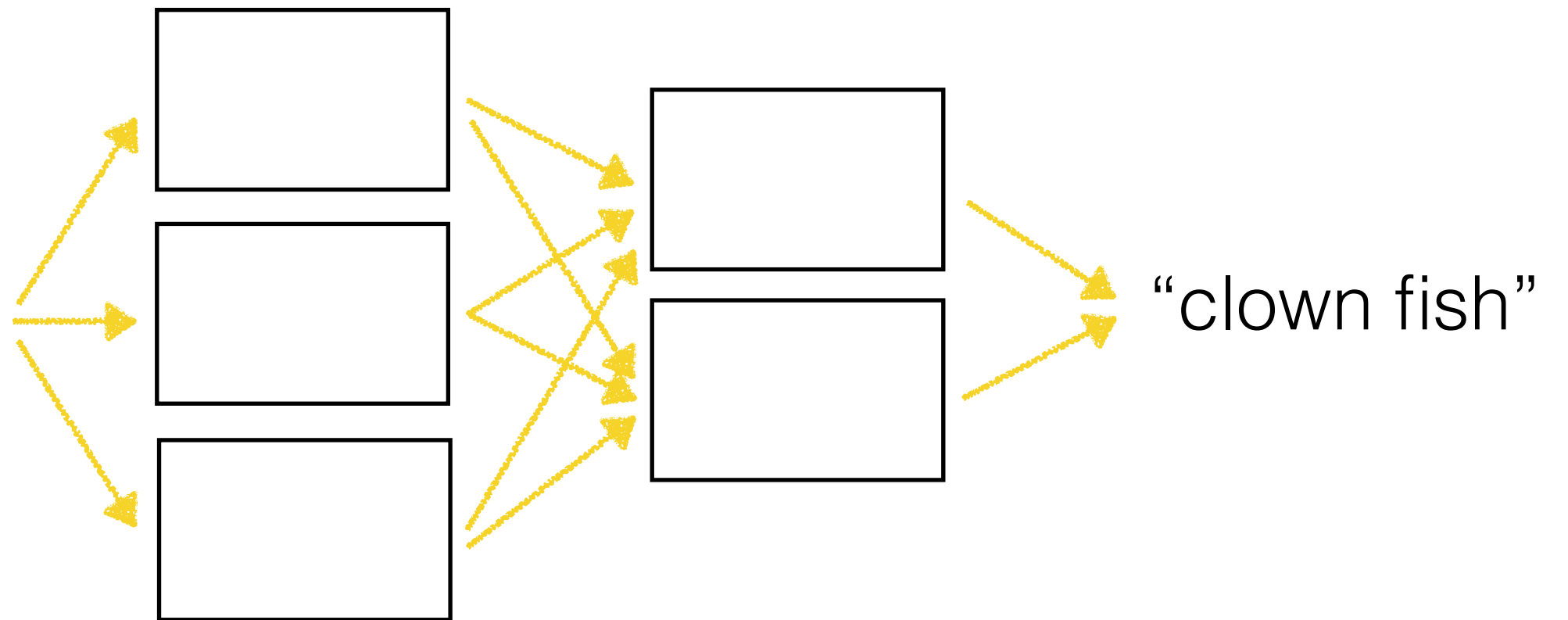
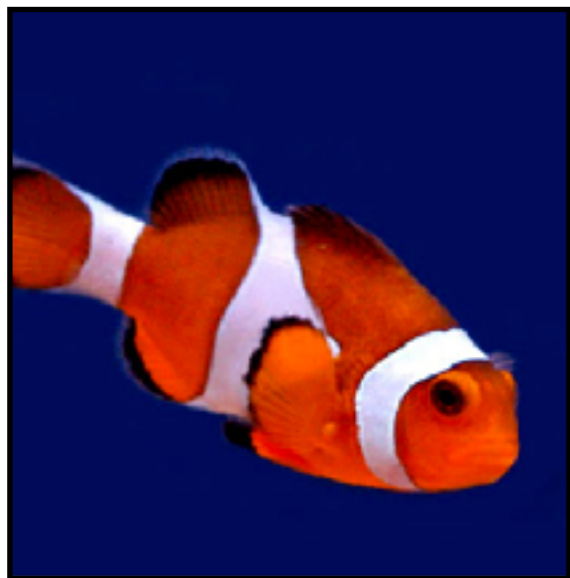
Object Recognition

Learned



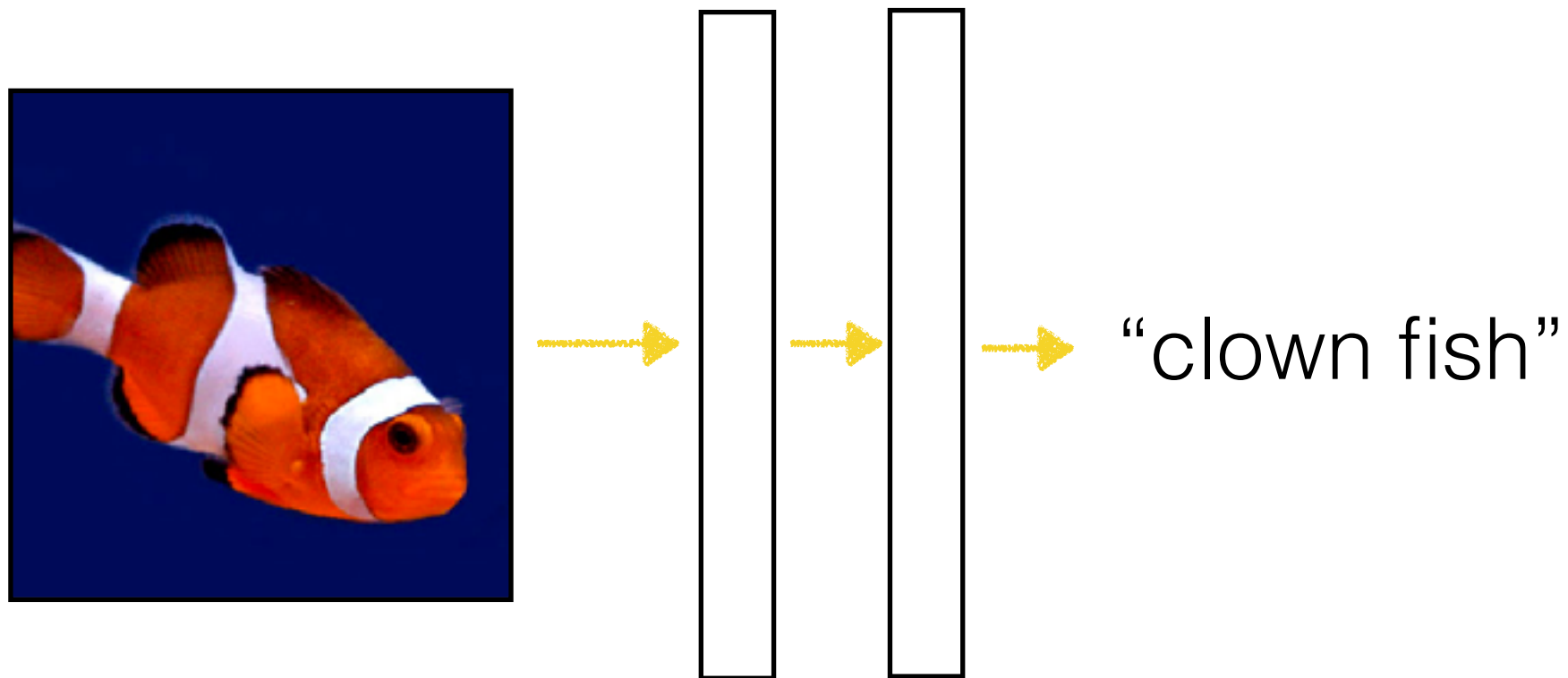
Neural Network

Learned



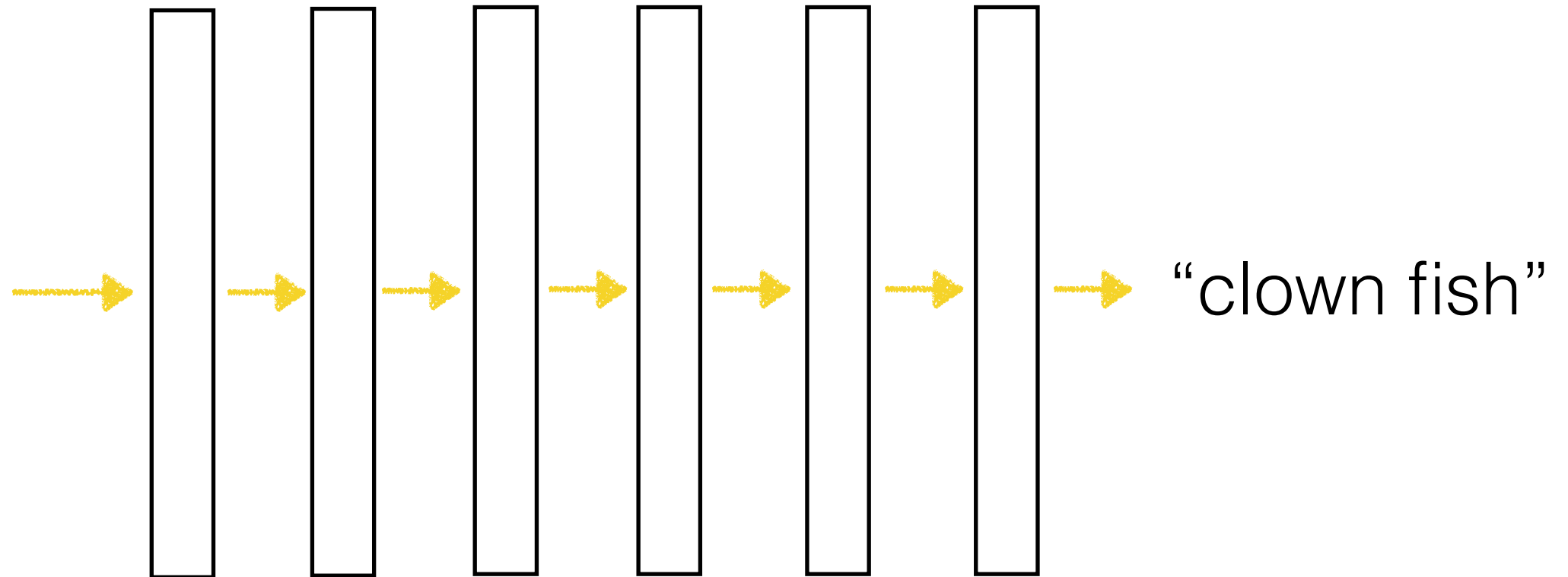
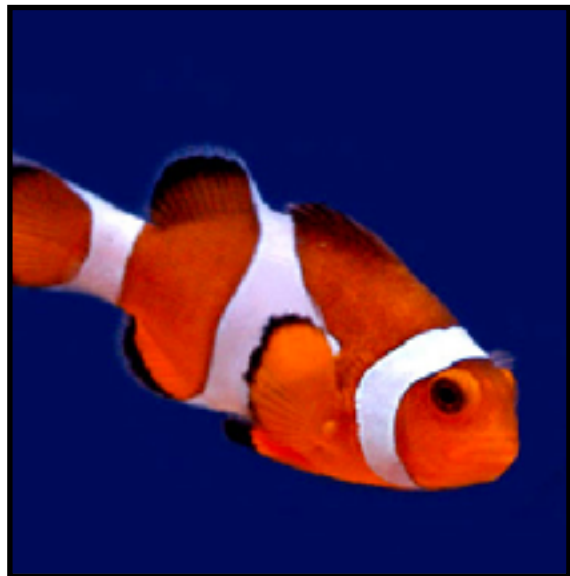
Neural Network

Learned

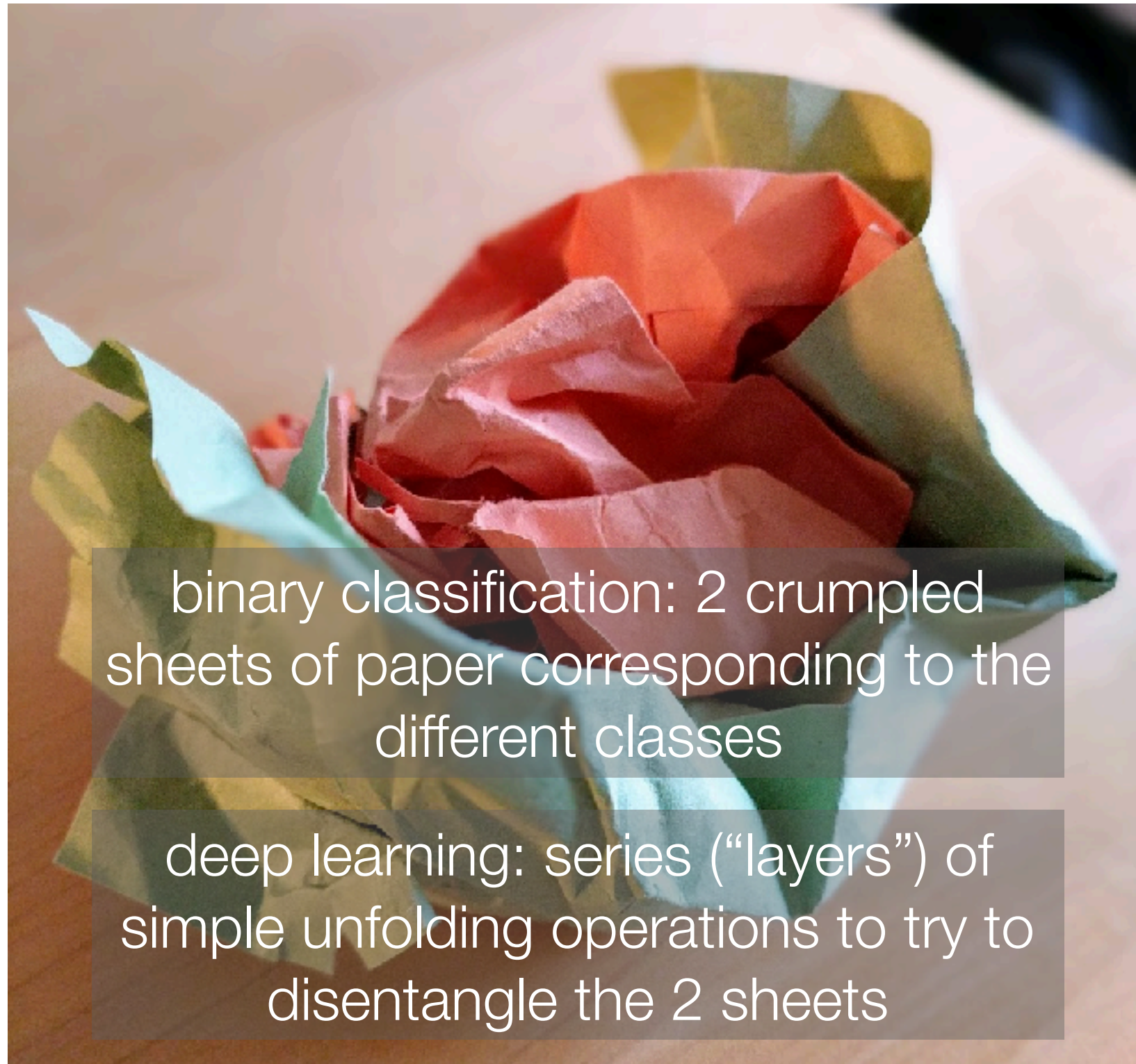


Deep Neural Network

Learned



Crumpled Paper Analogy



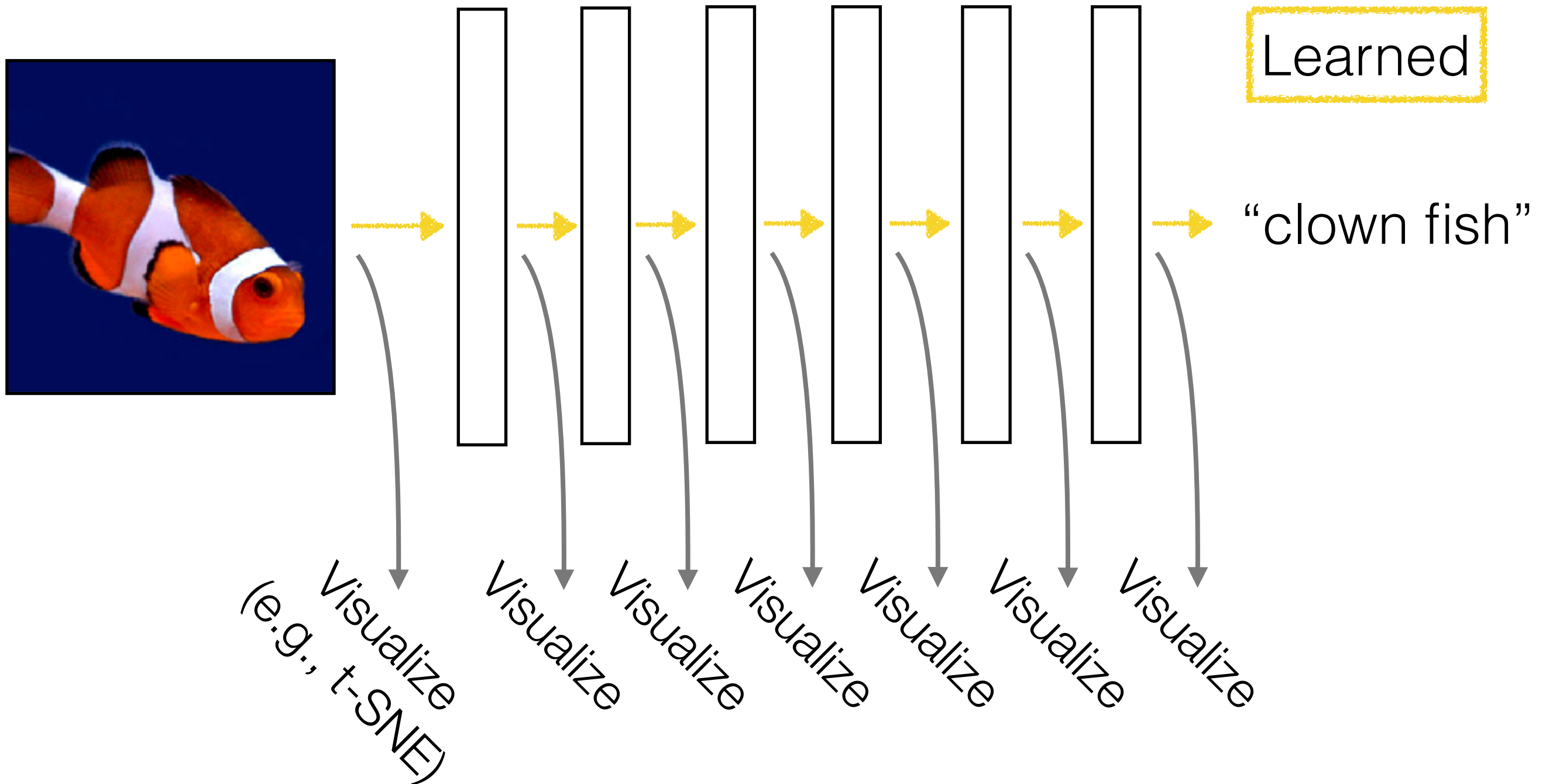
binary classification: 2 crumpled sheets of paper corresponding to the different classes

deep learning: series (“layers”) of simple unfolding operations to try to disentangle the 2 sheets

Analogy: Francois Chollet, photo: George Chen

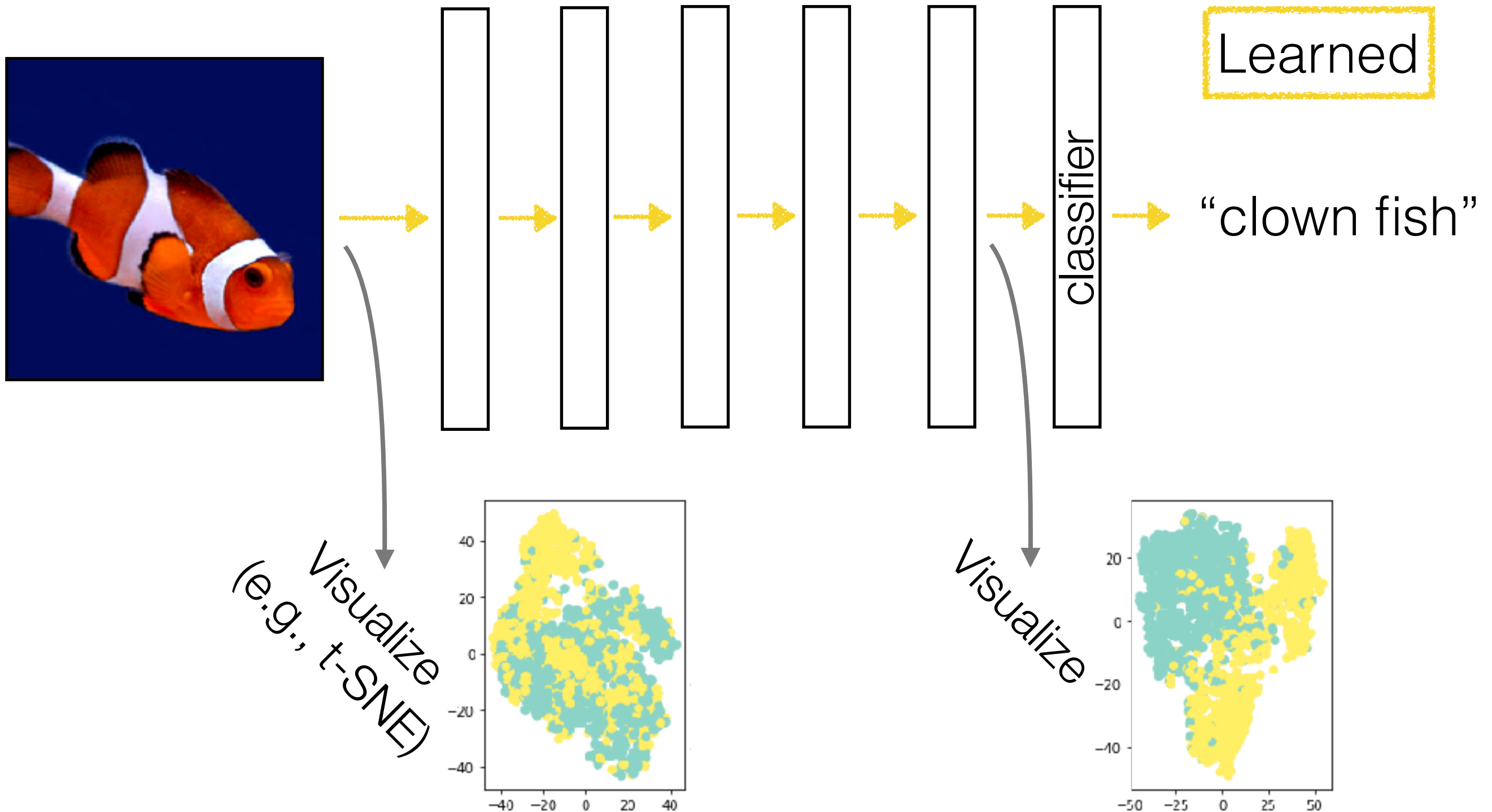
Representation Learning

Each layer's output is *another way we could represent the input data*



Representation Learning

Each layer's output is *another way we could represent the input data*



Why Does Deep Learning Work?

Actually the ideas behind deep learning are old (~1980's)

- Big data



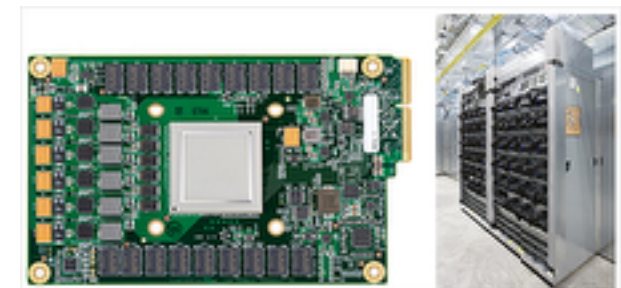
- Better hardware



CPU's
& Moore's law



GPU's



TPU's

- Better algorithms

Structure Present in Data Matters

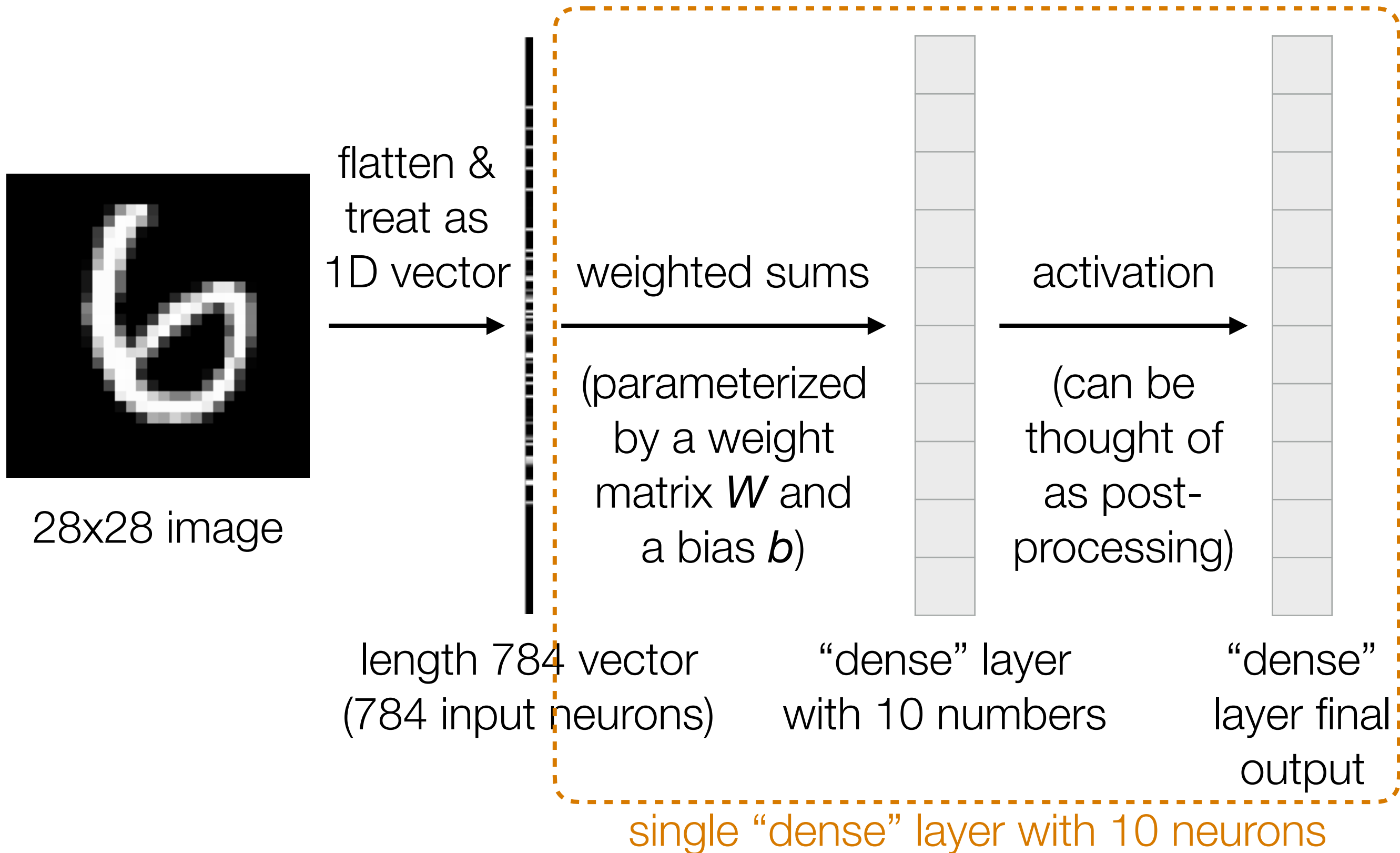
Neural nets aren't doing black magic

- **Image analysis:** convolutional neural networks (convnets) neatly incorporates basic image processing structure
- **Time series analysis:** recurrent neural networks (RNNs) incorporates ability to remember and forget things over time
 - Note: text is a time series
 - Note: video is a time series

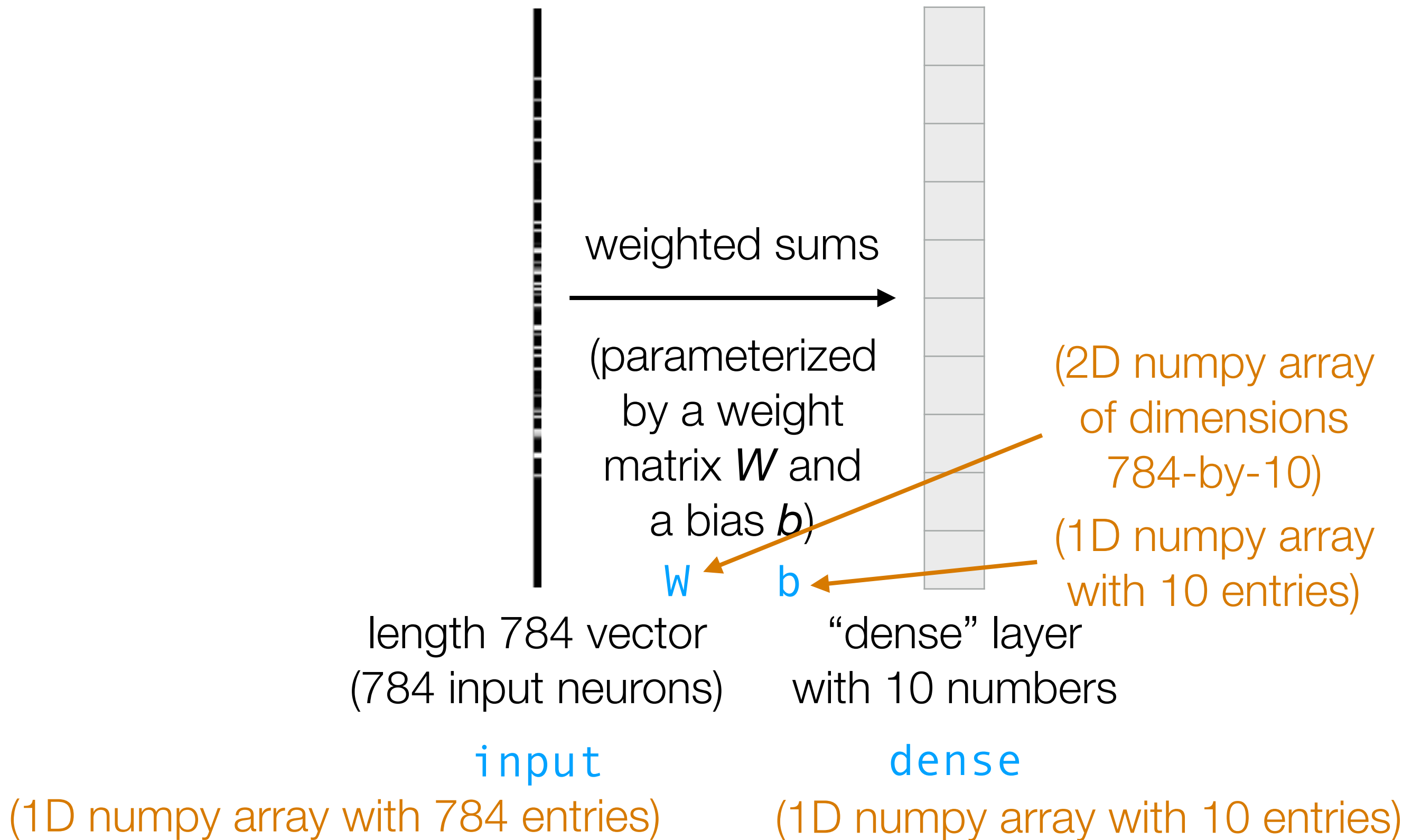
Handwritten Digit Recognition Example

Walkthrough of building a 1-layer and then a 2-layer neural net

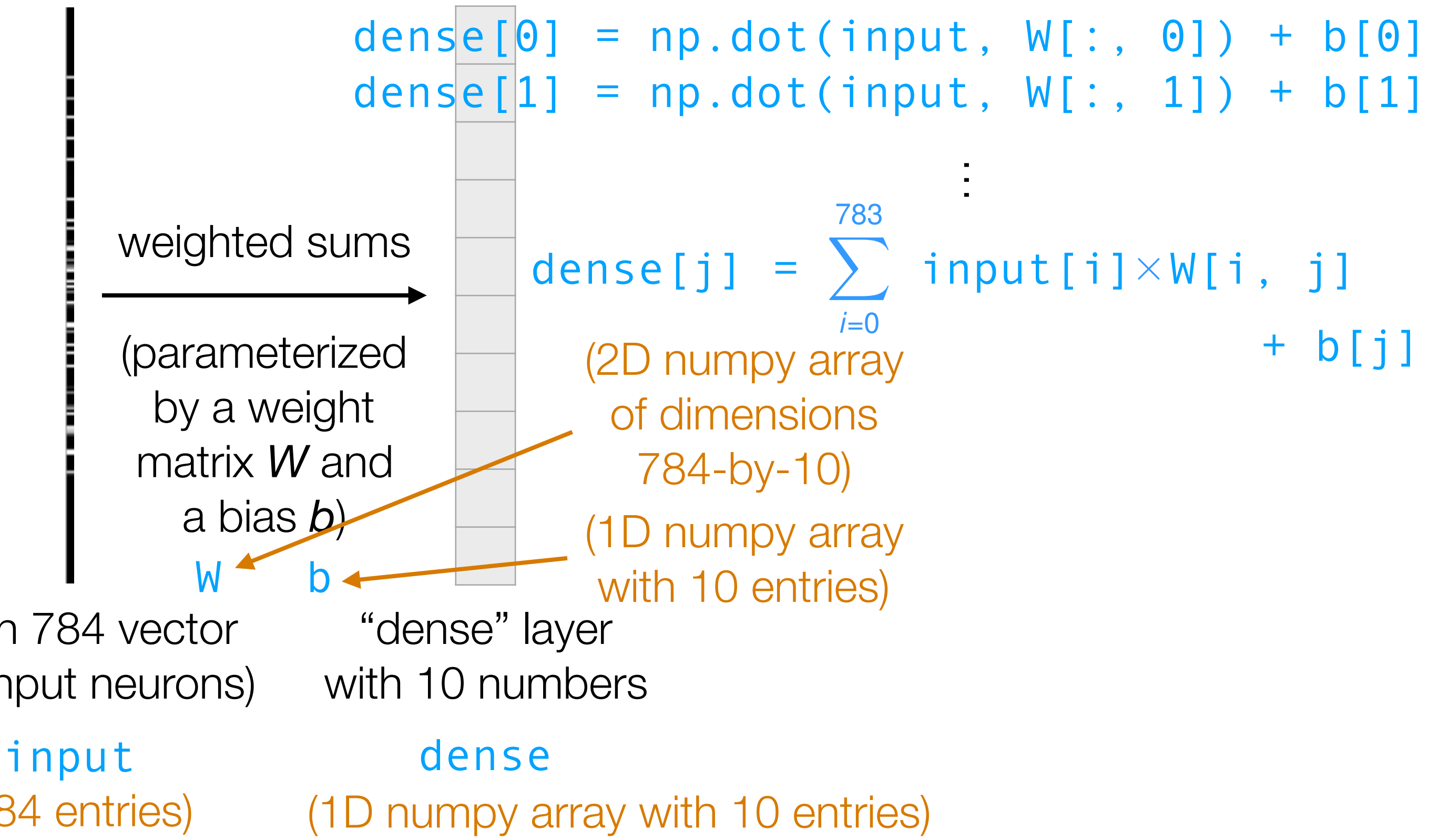
Handwritten Digit Recognition



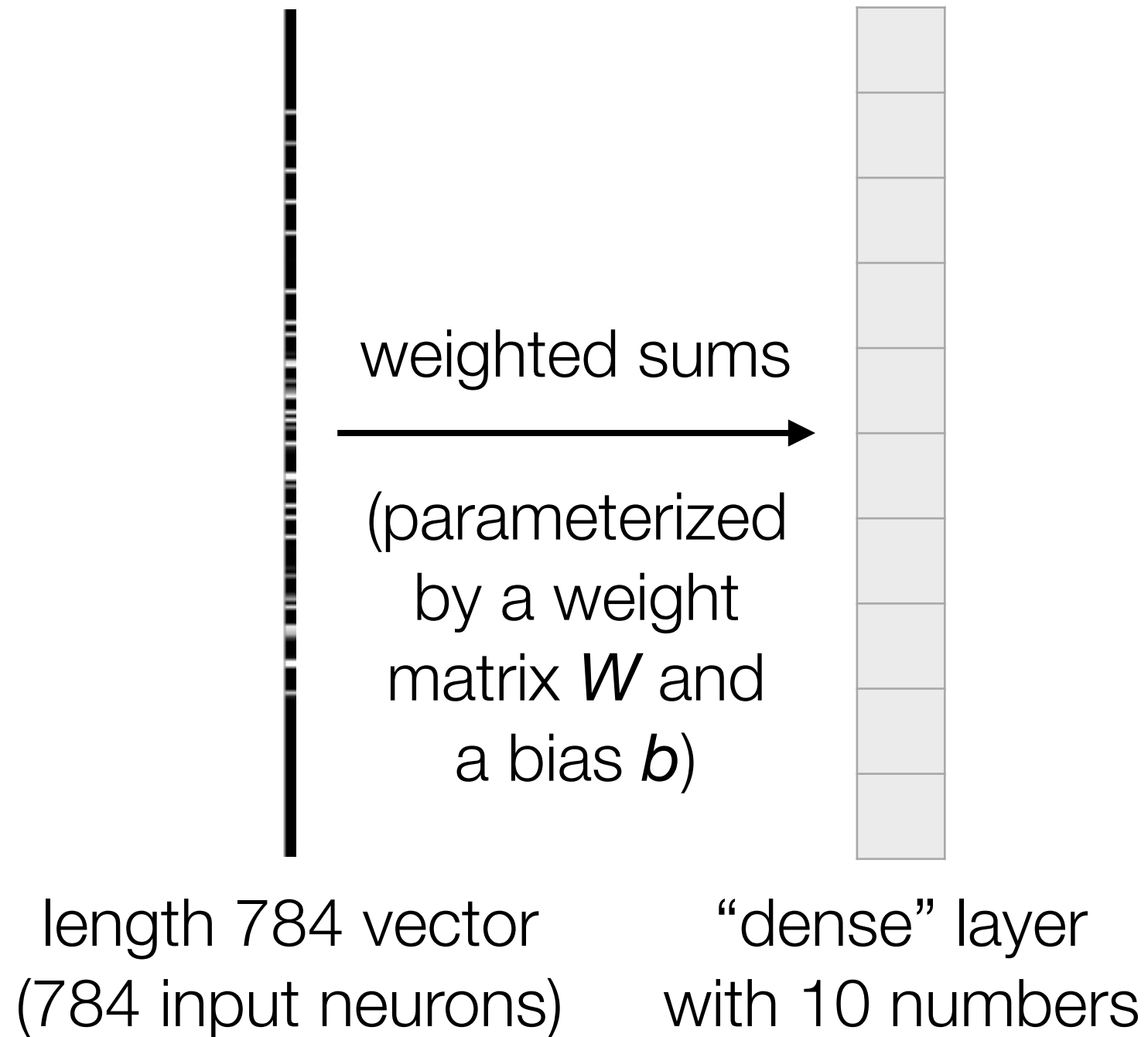
Handwritten Digit Recognition



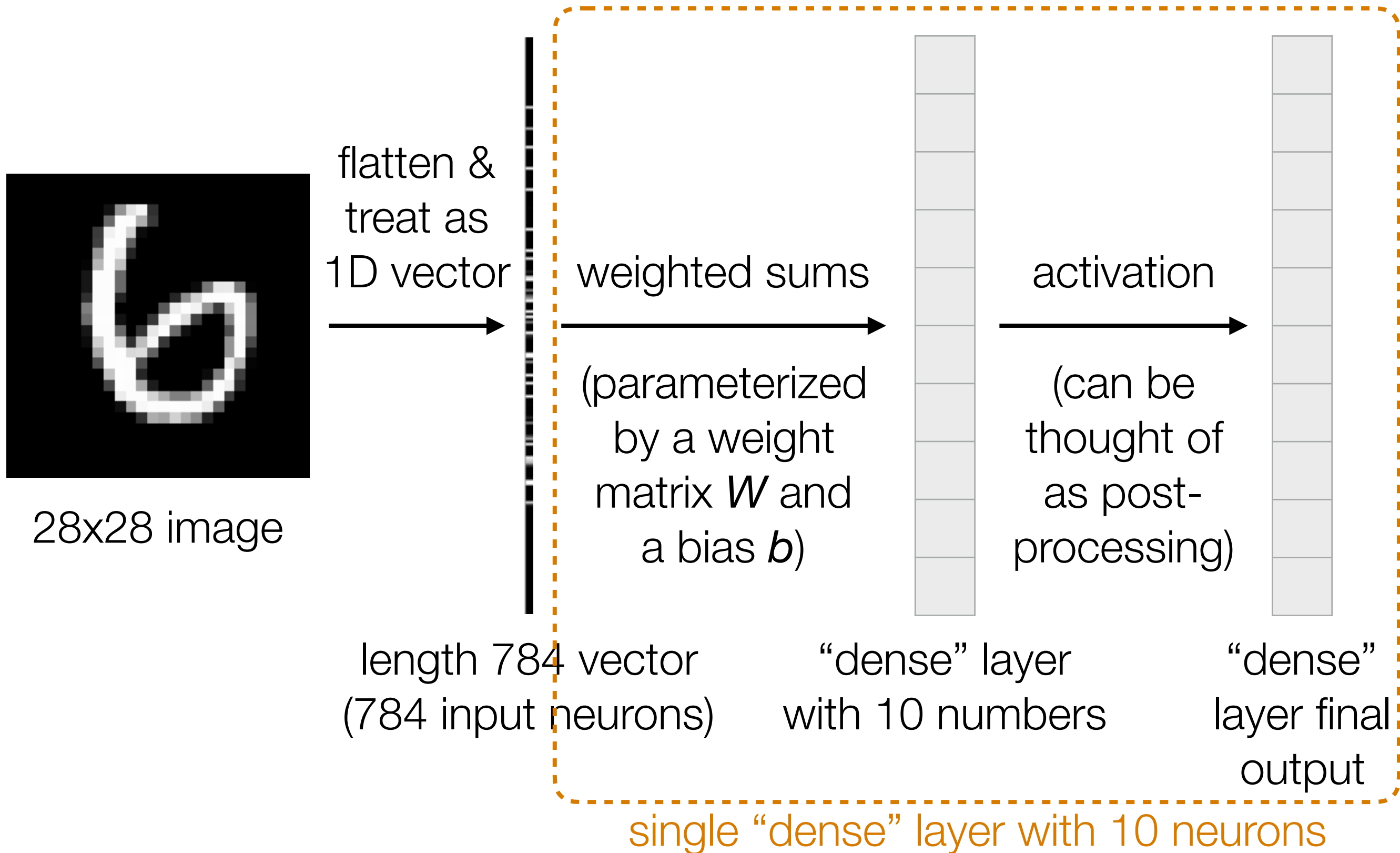
Handwritten Digit Recognition



Handwritten Digit Recognition



Handwritten Digit Recognition



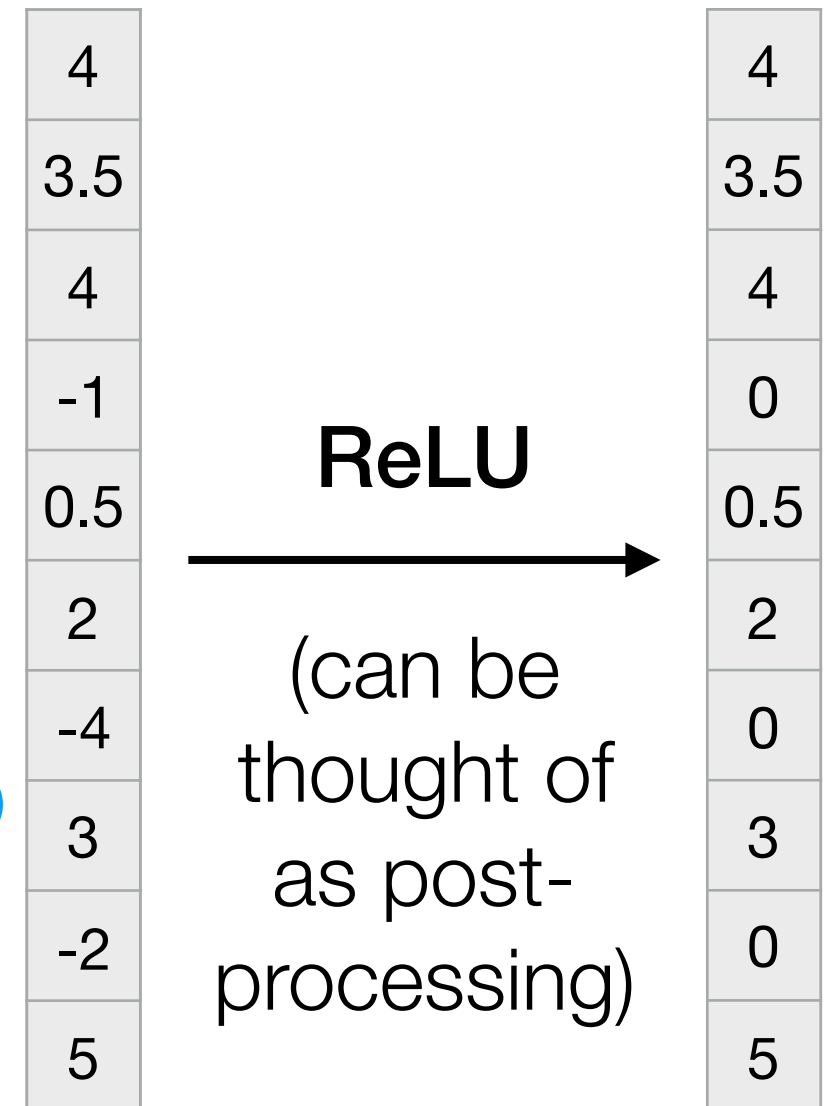
Handwritten Digit Recognition

Many different activation functions possible

Example: **Rectified linear unit (ReLU)**

zeros out entries that are negative

```
dense_final = np.maximum(0, dense)
```



“dense” layer
with 10 numbers

`dense`

“dense”
layer final
output

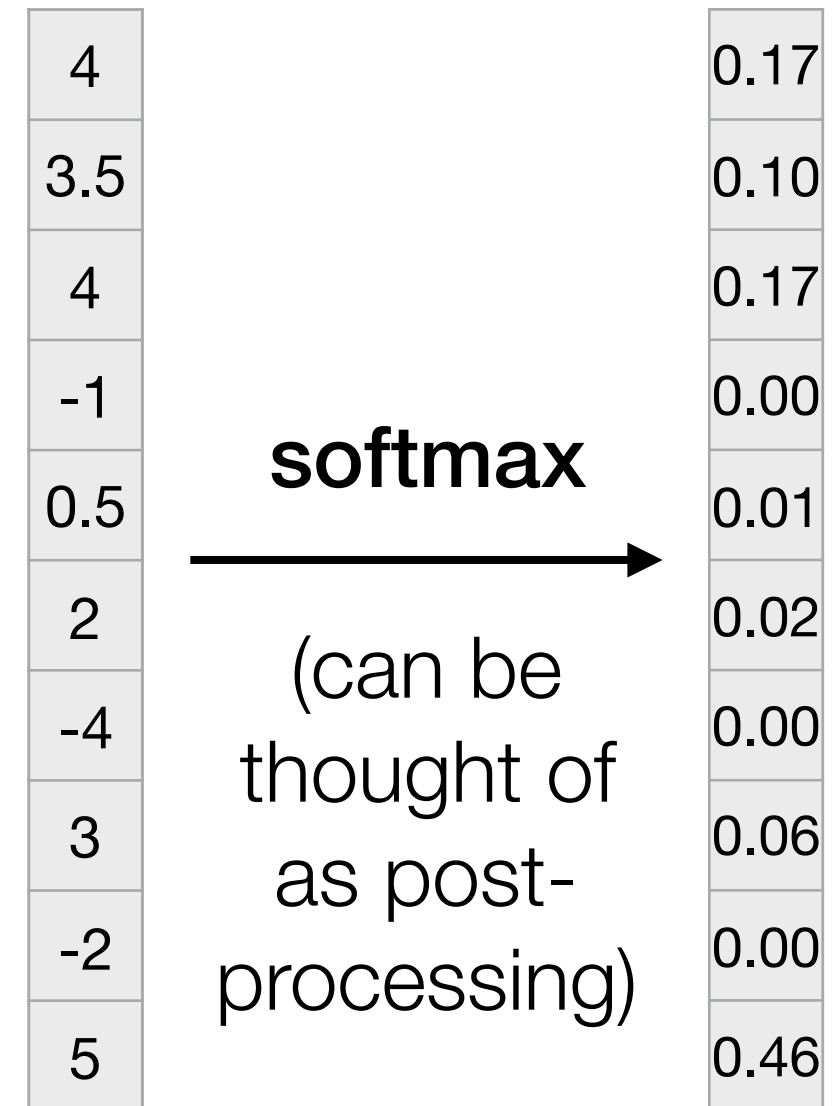
`dense_final`

Handwritten Digit Recognition

Many different activation functions possible

Example: **softmax** turns the entries in the dense layer (prior to activation) into a probability distribution (using the “softmax” transformation)

```
dense_exp = np.exp(dense)
dense_exp /= np.sum(dense_exp)
dense_final = dense_exp
```



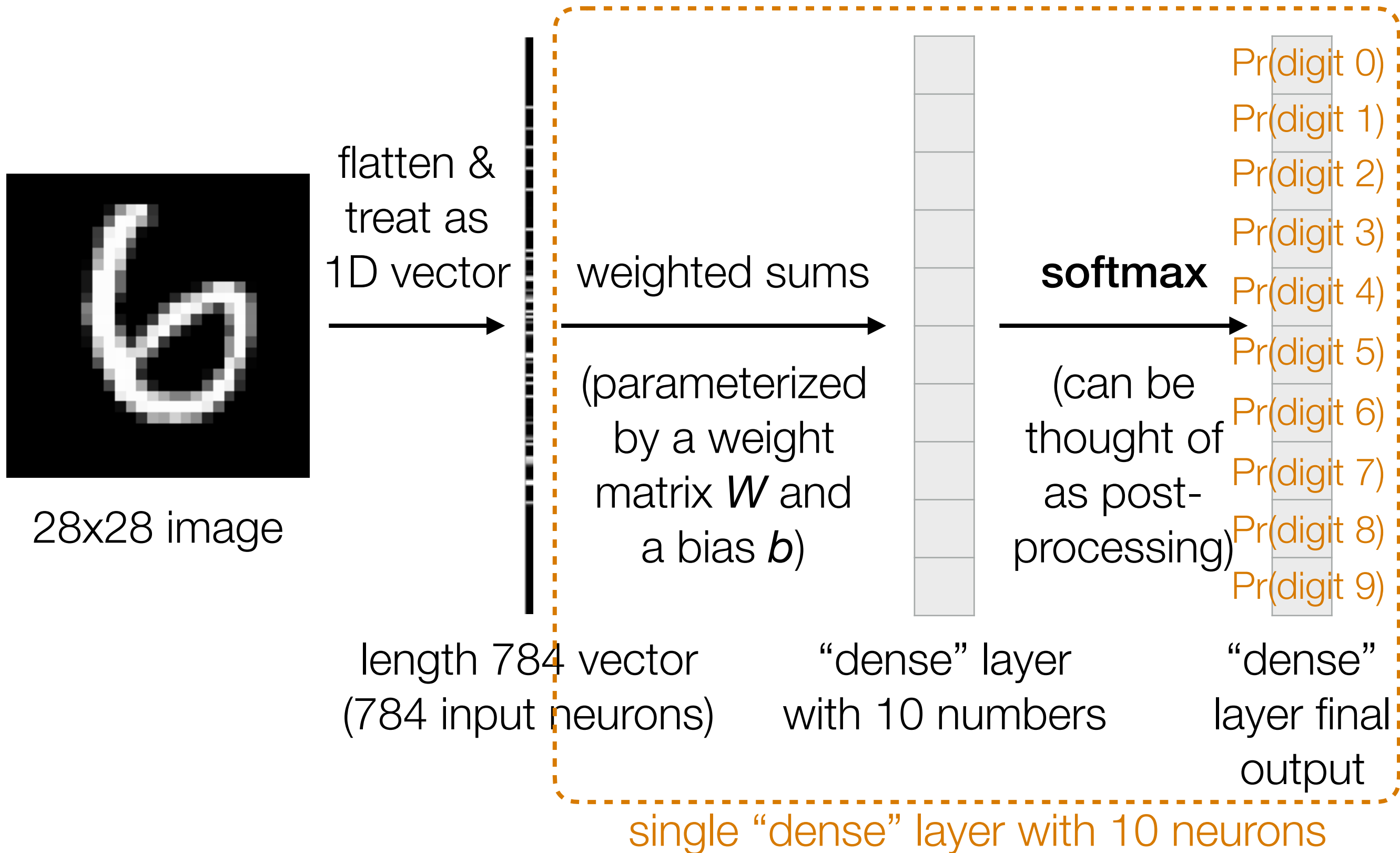
“dense” layer
with 10 numbers

`dense`

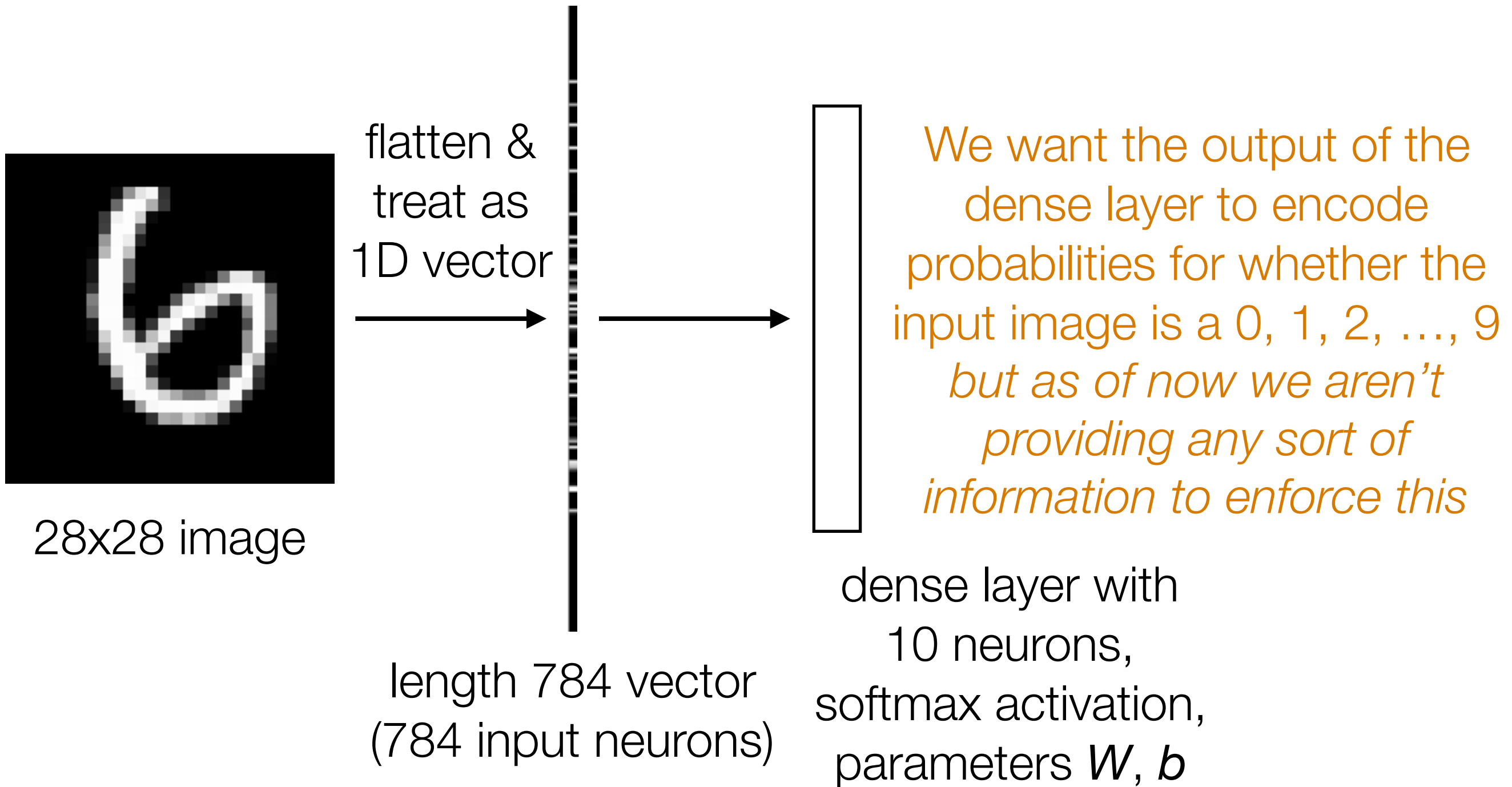
“dense”
layer final
output

`dense_final`

Handwritten Digit Recognition



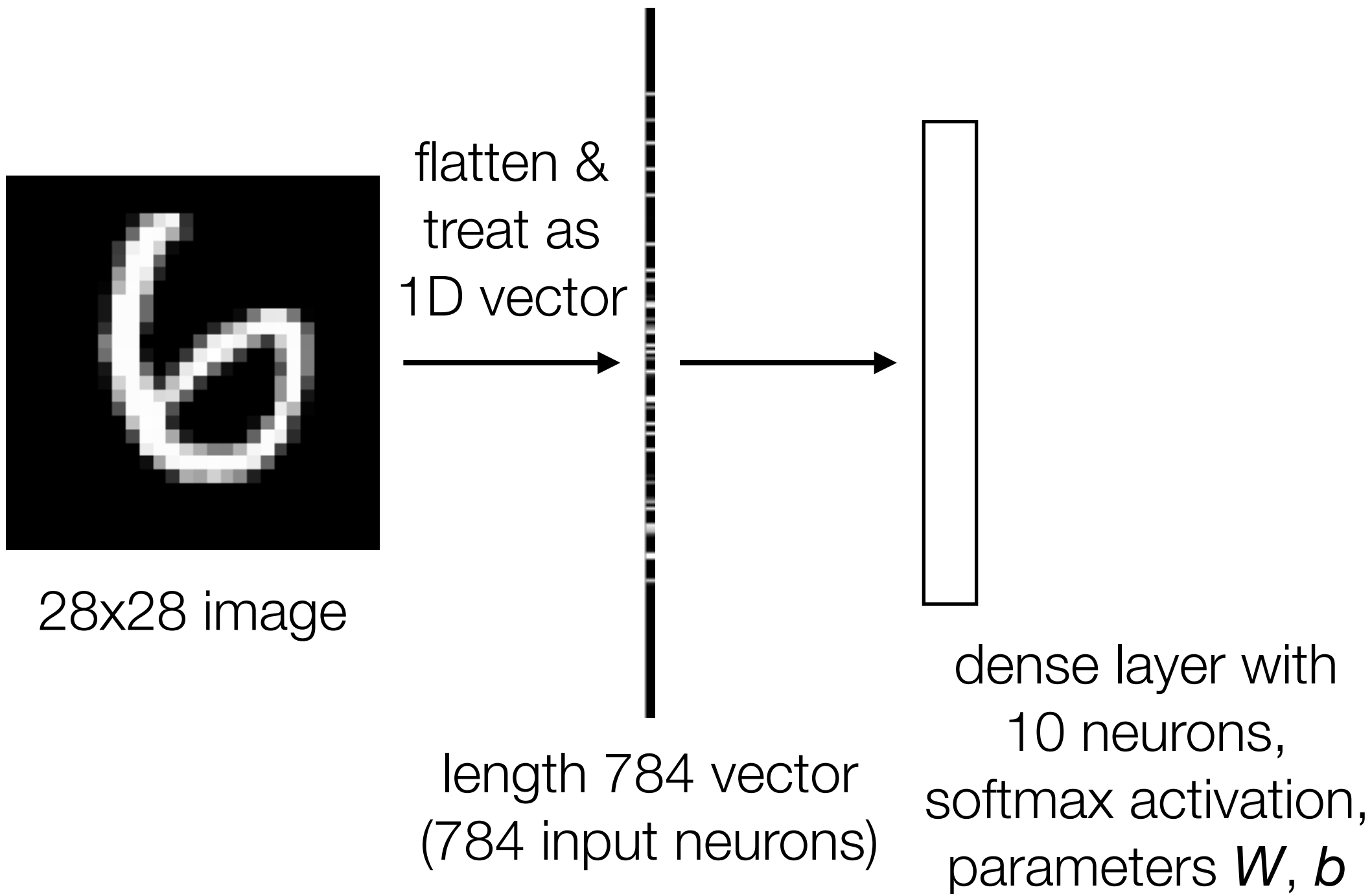
Handwritten Digit Recognition



Handwritten Digit Recognition

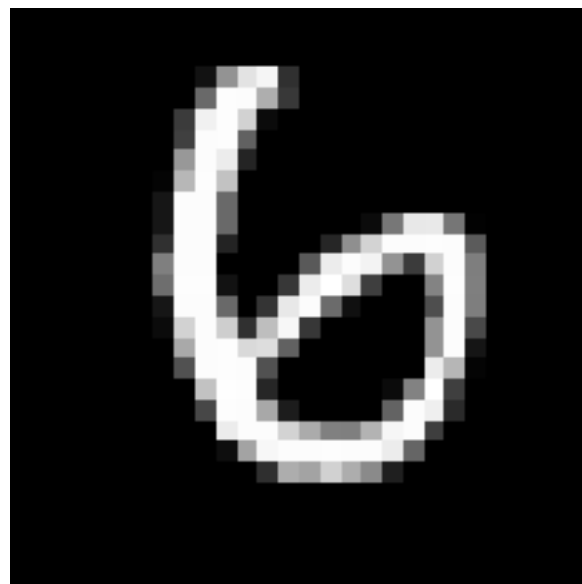
Demo part 1

Handwritten Digit Recognition



Handwritten Digit Recognition

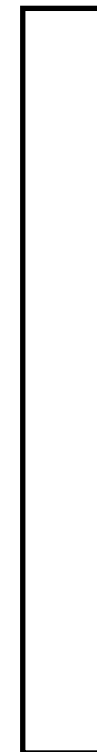
Training label: 6



28x28 image

Learning this neural net means learning W and b

flatten & treat as 1D vector



Loss/"error"



Error is averaged across training examples

length 784 vector (784 input neurons)

dense layer with 10 neurons, softmax activation, parameters W, b

Popular loss function for classification (> 2 classes): **categorical cross entropy**

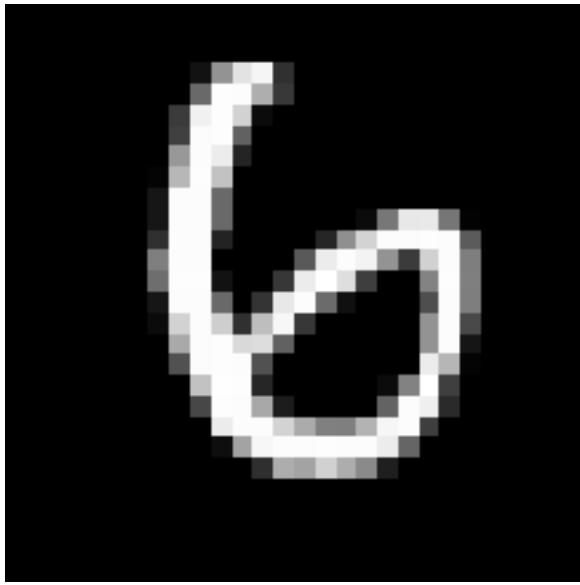
$$\log \frac{1}{\text{Pr}(\text{digit } 6)}$$

Handwritten Digit Recognition

Demo part 2

Handwritten Digit Recognition

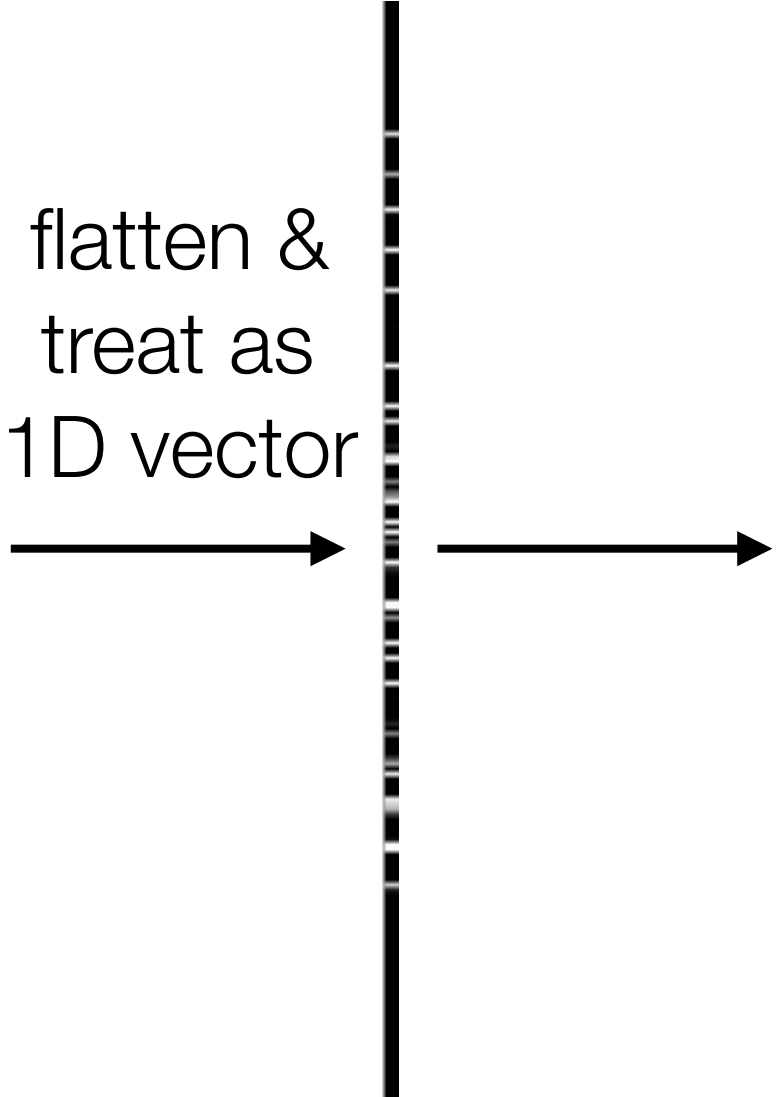
Training label: 6



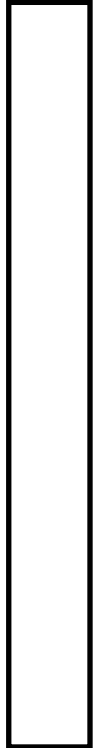
28x28 image

Learning this neural net means learning W and b

flatten & treat as 1D vector



length 784 vector (784 input neurons)



dense layer with 10 neurons, softmax activation, parameters W, b

Loss/"error"

Popular loss function for classification (> 2 classes): **categorical cross entropy**

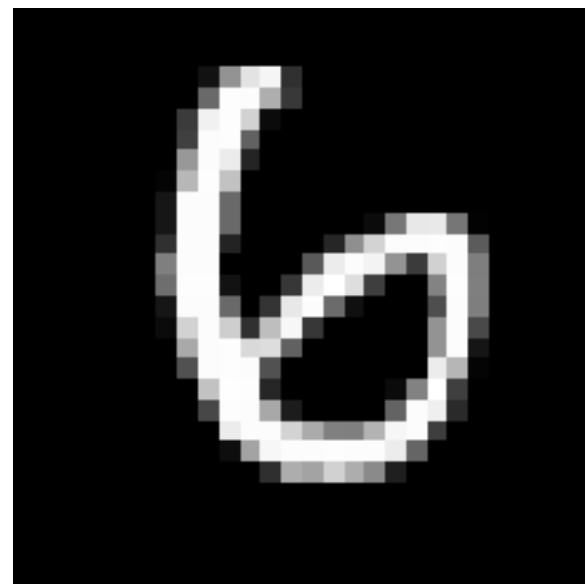
$$\log \frac{1}{\text{Pr}(\text{digit } 6)}$$

Error is averaged across training examples

error

Handwritten Digit Recognition

Training label: 6



28x28 image

length 784 vector
(784 input neurons)

Learning this neural net means learning parameters of both dense layers!



dense layer with 512 neurons, ReLU activation

dense layer with 10 neurons, softmax activation

Loss/"error"

Popular loss function for classification (> 2 classes): **categorical cross entropy**

$$\log \frac{1}{\text{Pr}(\text{digit } 6)}$$

Error is averaged across training examples

error

Handwritten Digit Recognition

Demo part 3