

**End of Days**

# What was this course about?

- Machine Learning
- Deep Learning
- Really just thinking, the process

## Along the way we

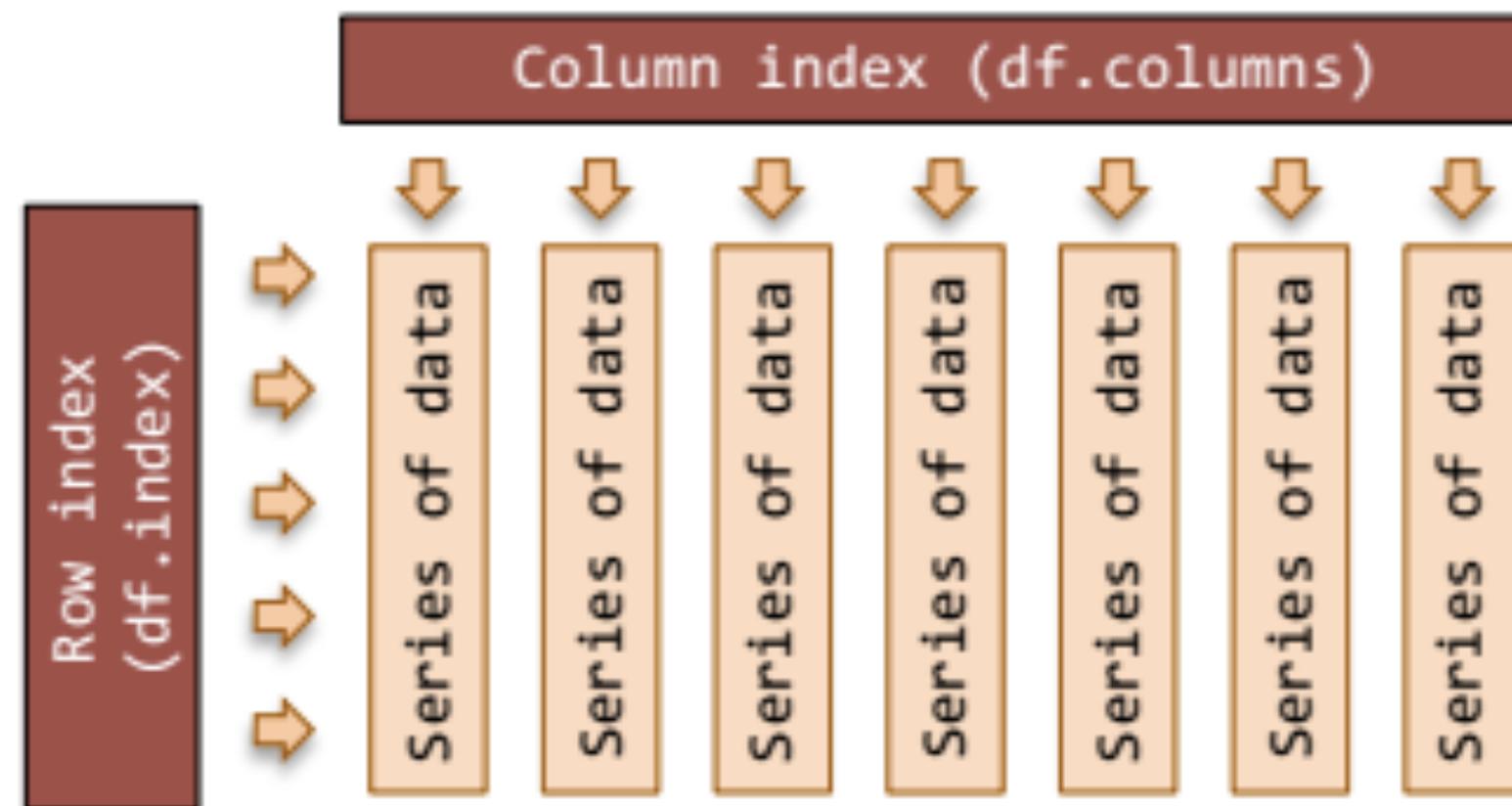
- learn how to create and regularize models
- learn how to optimize objective functions such as loss functions using Stochastic Gradient Descent
- learnt how to create (simple) architectures to solve problems
- learnt how to transfer from one model to the next

# Concepts running through:

- Fitting parameters vs hyperparameters
- Regularization
- Validation testing
- Stochastic Gradient Descent
- Learning representations

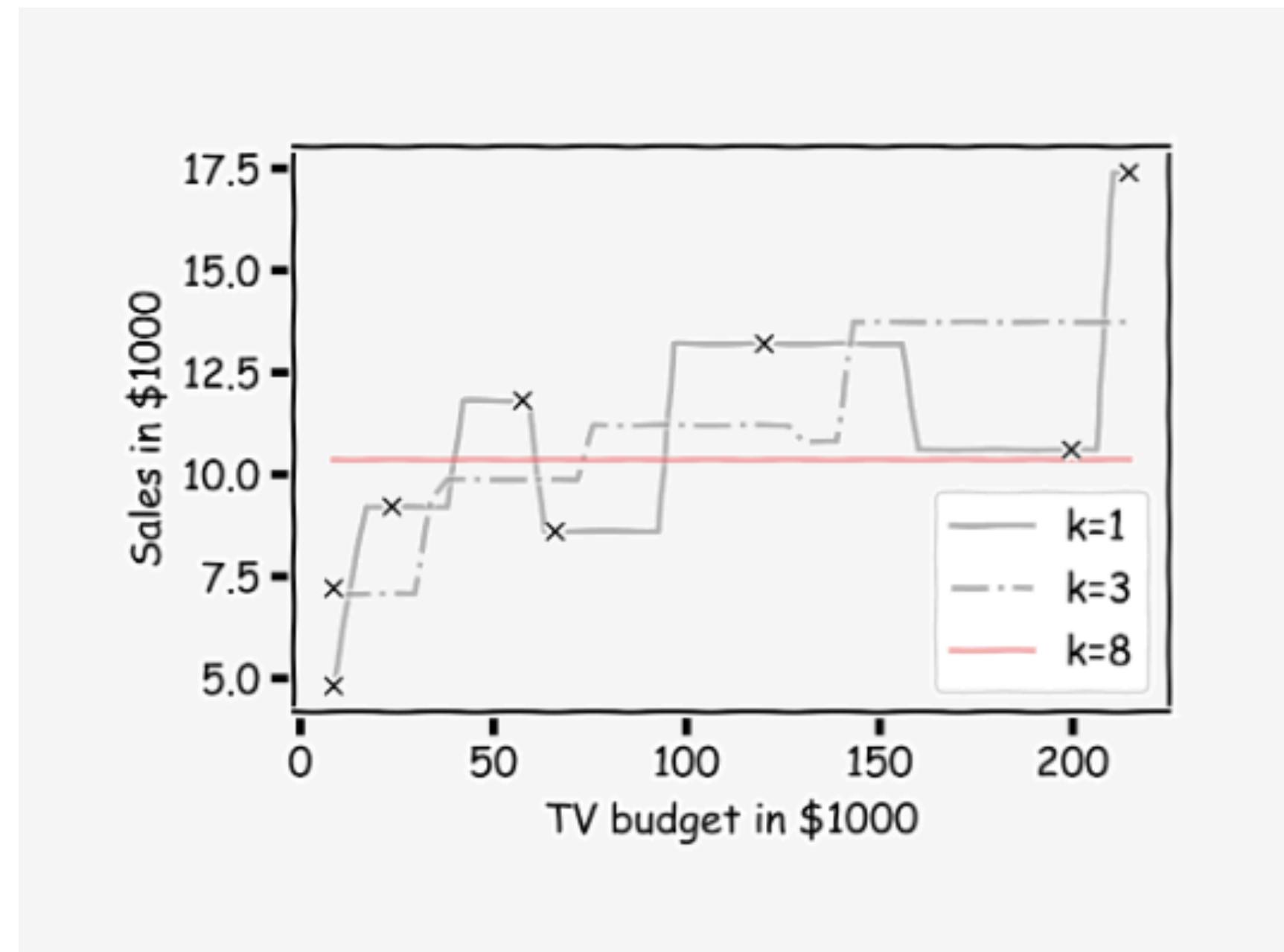
# Tablular data and Pandas

DataFrame object: a two-dimensional table of data with column and row indexes. The columns are made up of pandas Series objects.

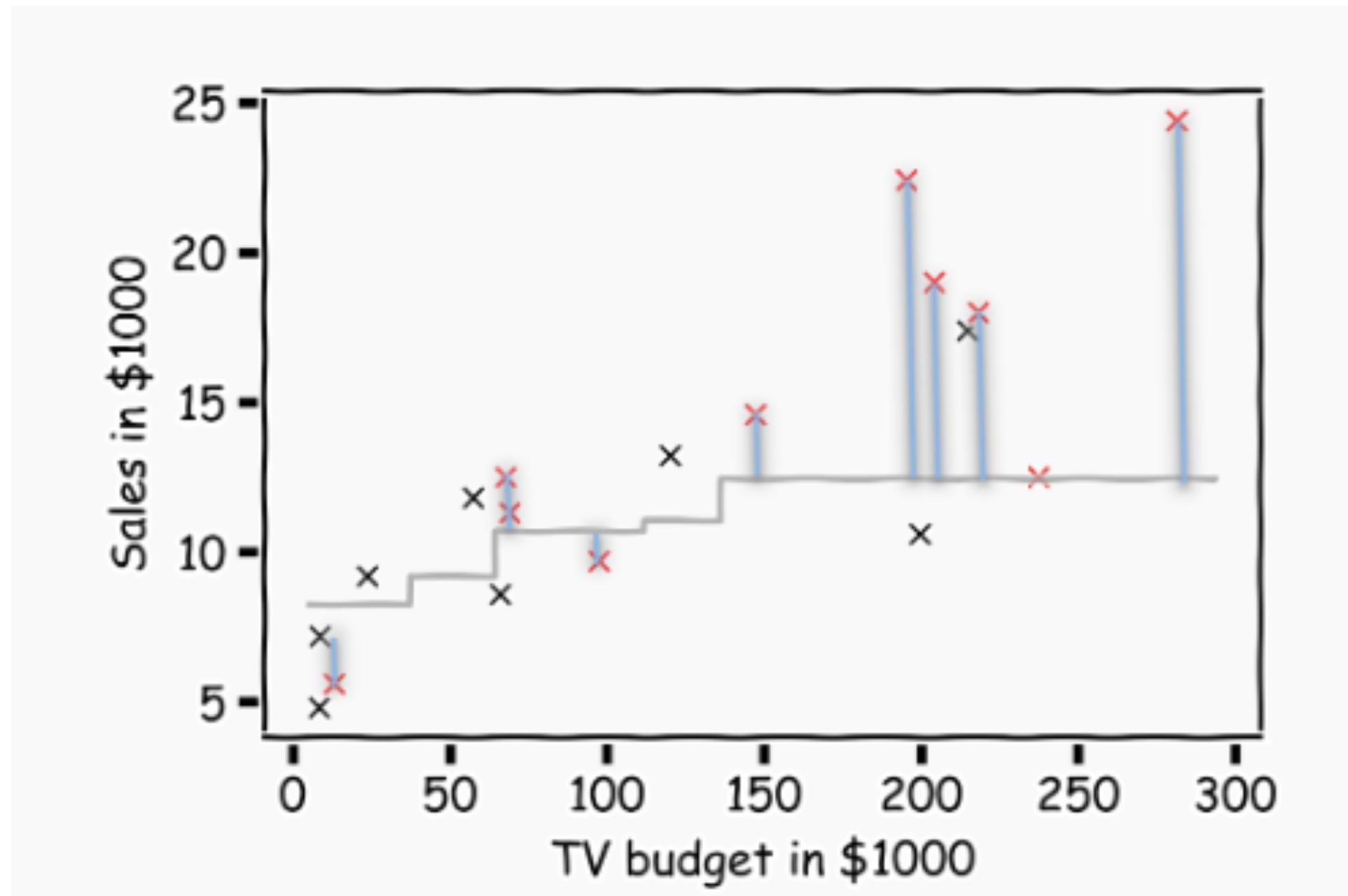


# Regression

# KNN Regression



# Residuals and their minimization



## How to fit: sklearn

```
Xtrain, Xtest, ytrain, ytest = train_test_split(X, y, test_size=0.2)
lr2 = LinearRegression().fit(Xtrain, ytrain)
r2_test = r2_score(ytest, lr.predict(Xtest))
r2_train = r2_score(ytrain, lr.predict(Xtrain))
```

# How to fit: Keras

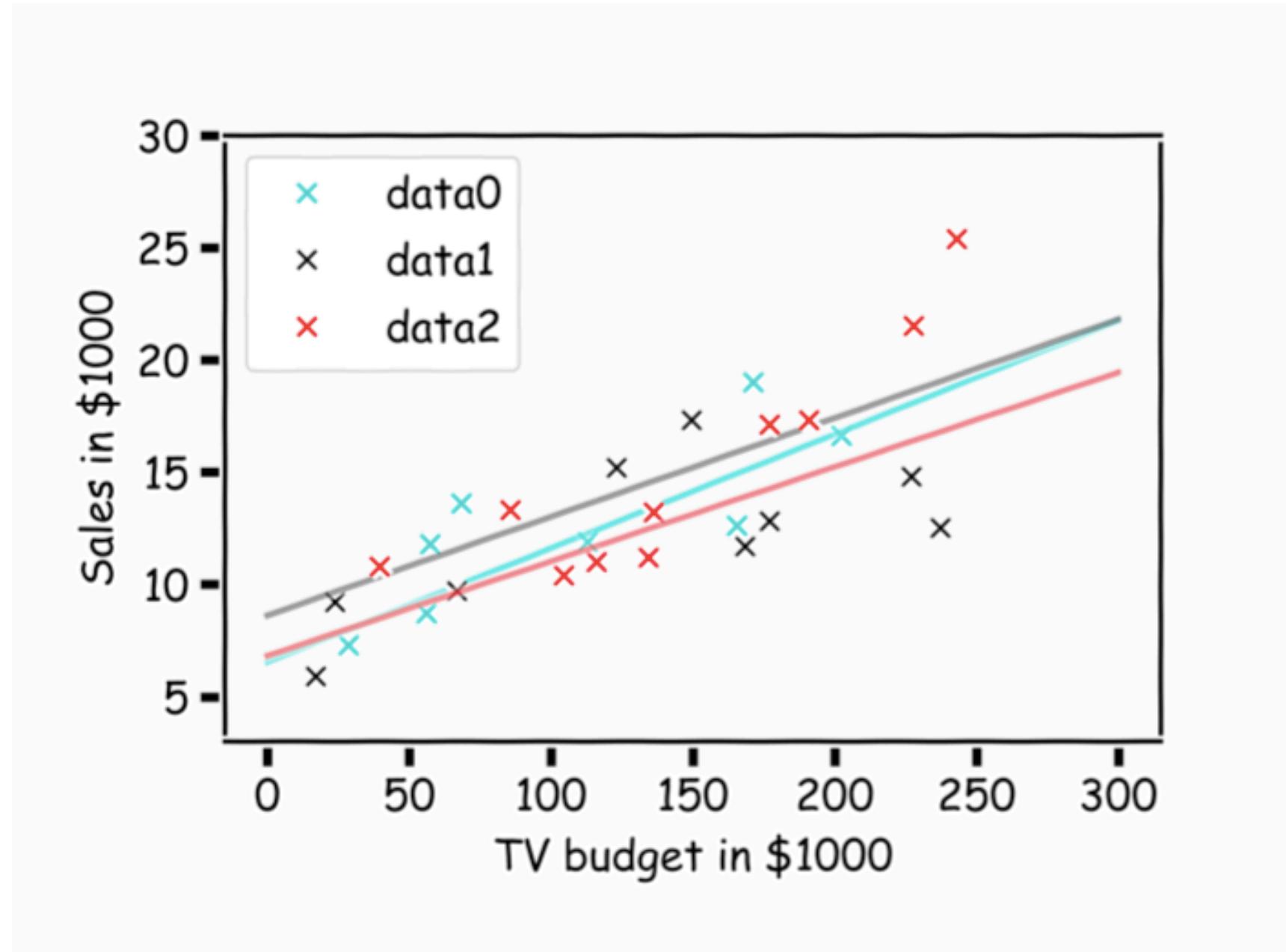
```
inputs_placeholder = Input(shape=(1,))  
outputs_placeholder = Dense(1, activation='linear')(inputs_placeholder)  
  
m = Model(inputs=inputs_placeholder, outputs=outputs_placeholder)  
m.compile(optimizer='sgd', loss='mean_squared_error', metrics=['mae', 'accuracy'])  
m.summary()
```

# Frequentist Statistics

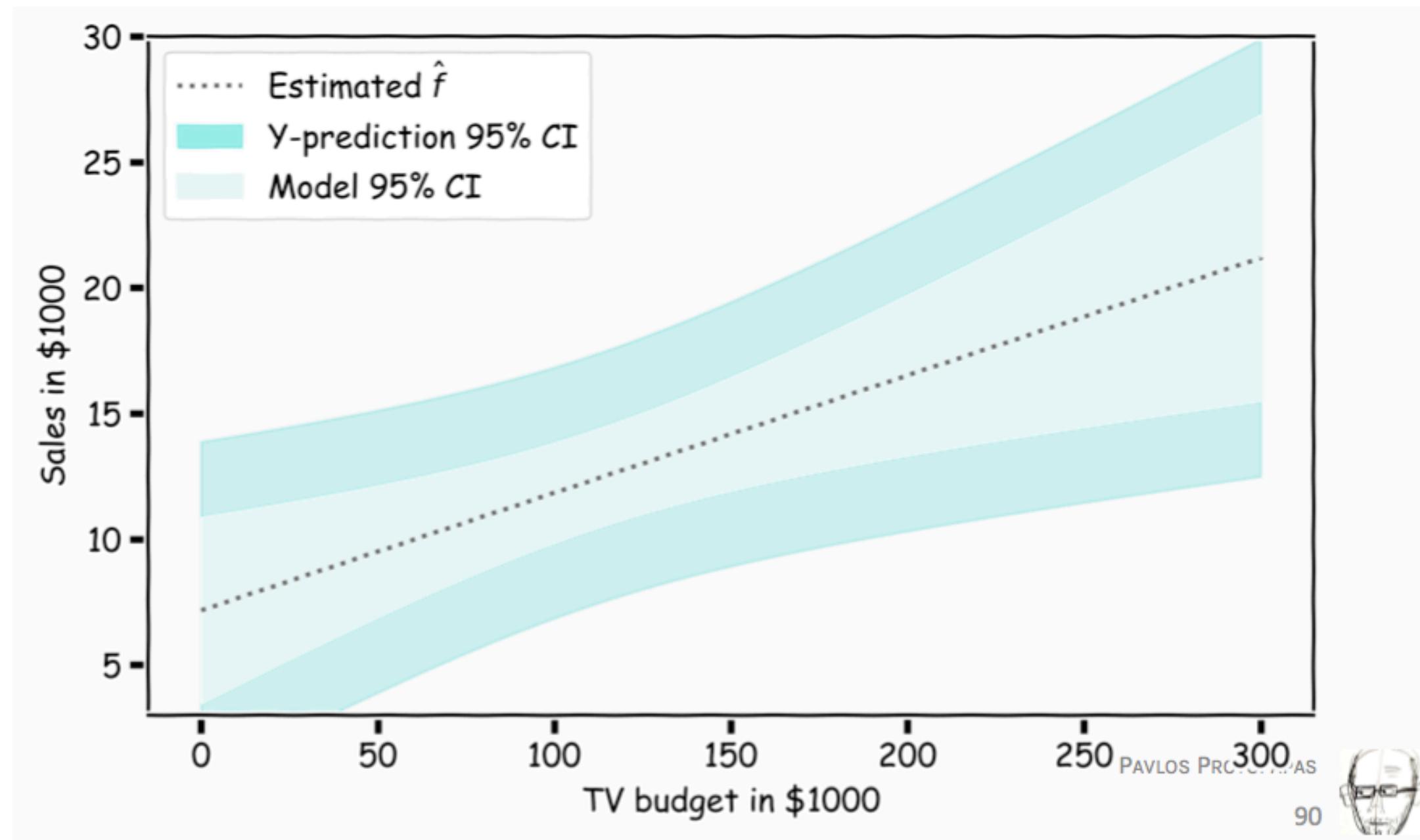
"data is a **sample** from an existing **population**"

- data is stochastic, variable; parameters fixed
- fit a parameter
- samples (or bootstrap) induce a sampling distribution on any estimator
- example of a very useful estimator: MLE

# Multiple fits from multiple samples

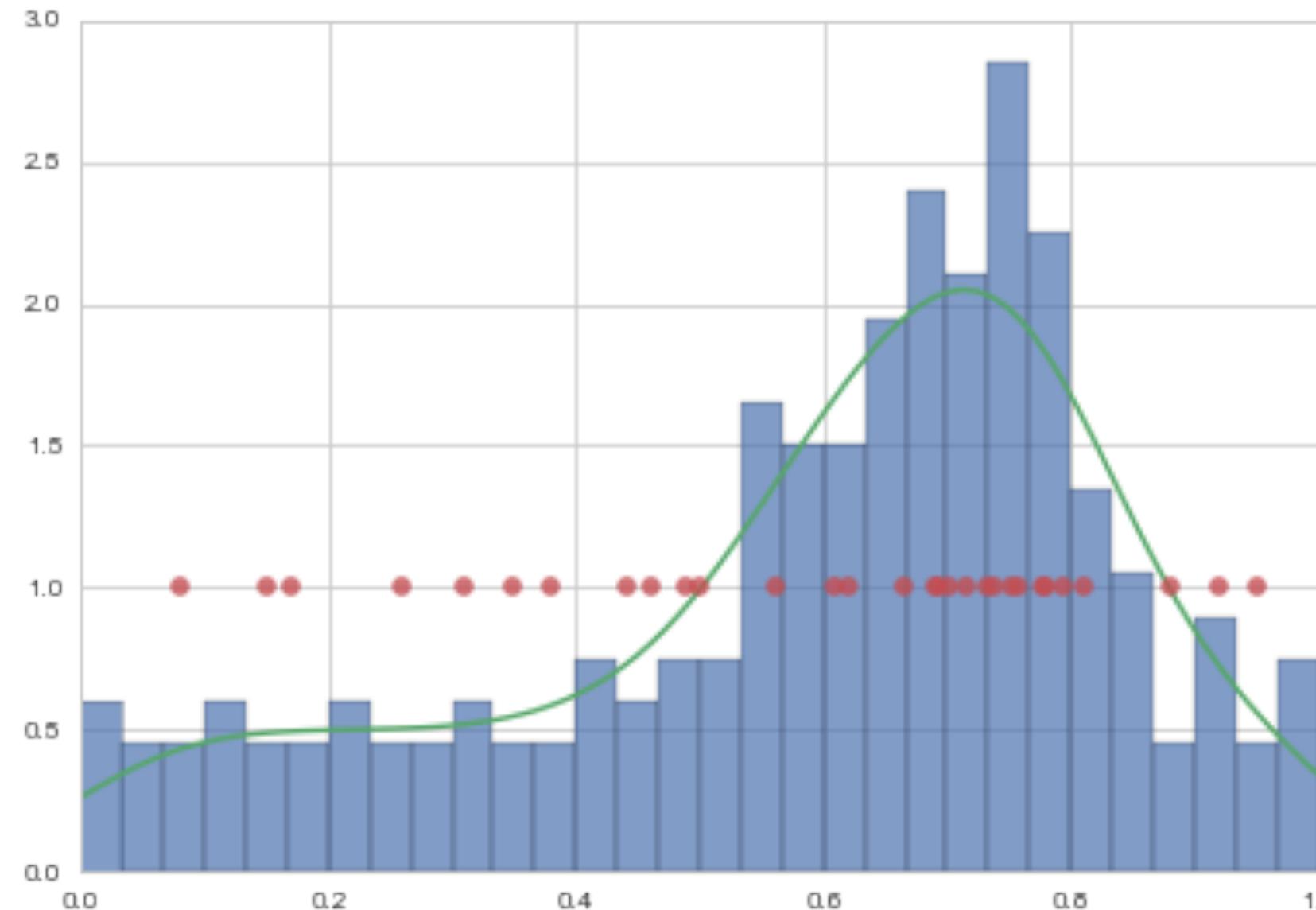


# Regression line uncertainty vs prediction uncertainty



learning

# Statement of the Learning Problem

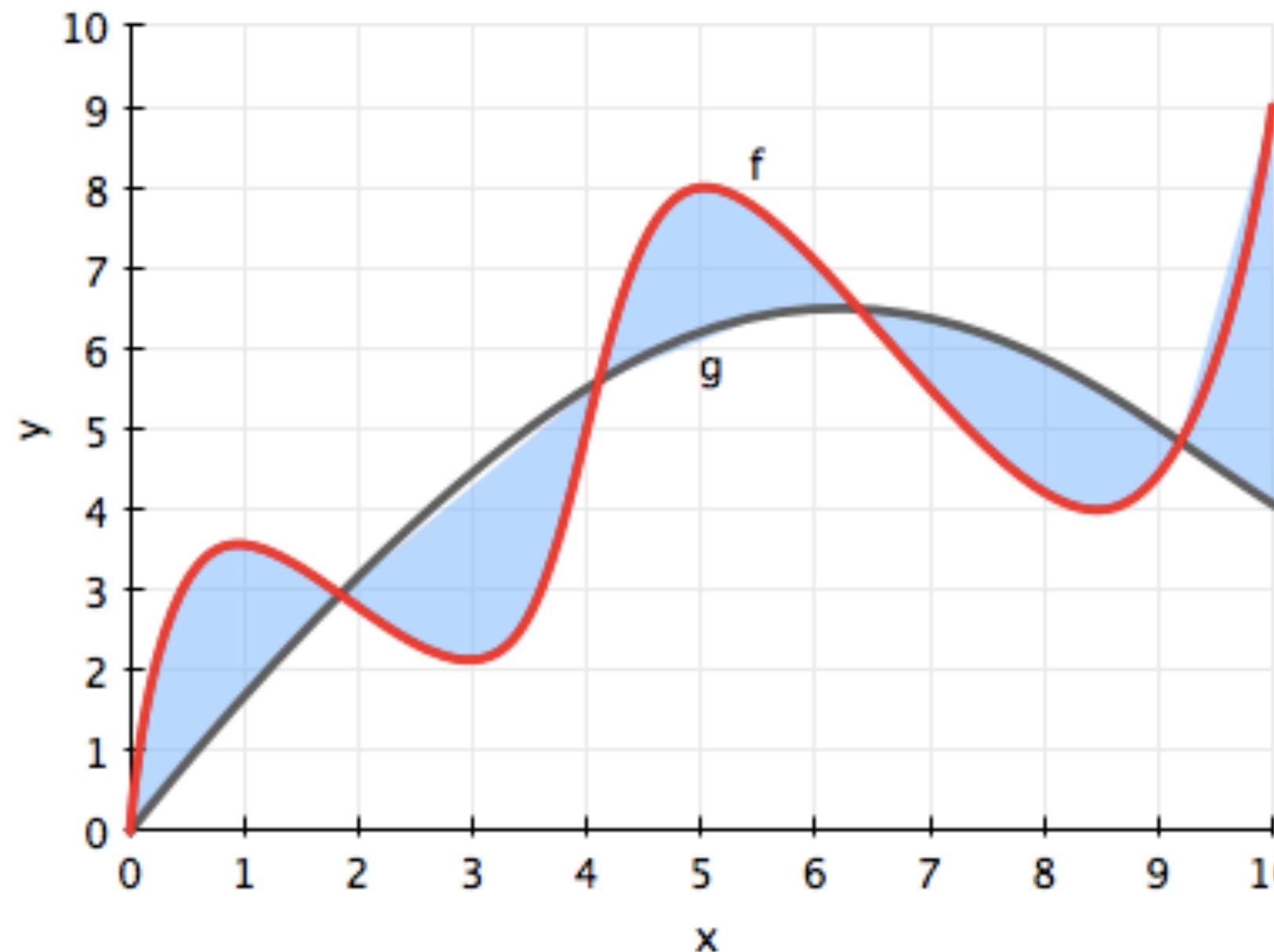


The sample must be representative of the population!

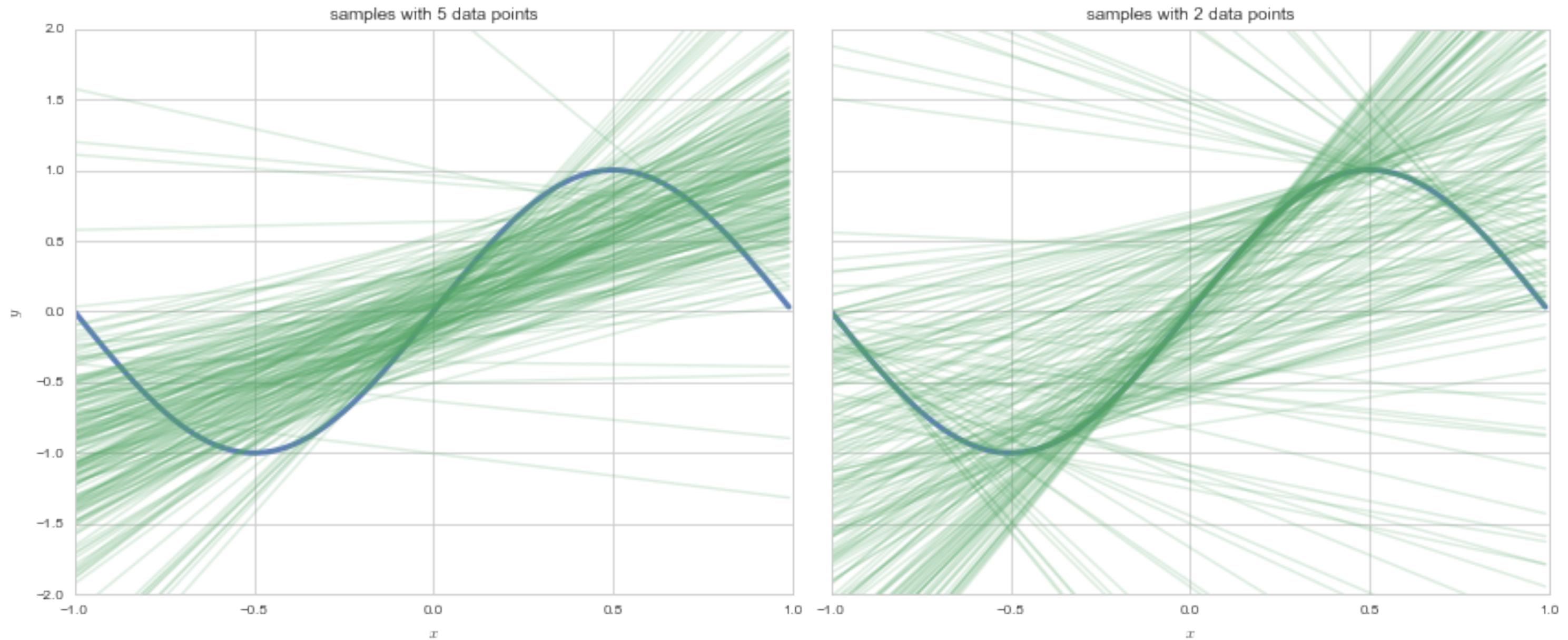
- A :  $R_{\mathcal{D}}(g)$  smallest on  $\mathcal{H}$*
- B :  $R_{out}(g) \approx R_{\mathcal{D}}(g)$*

A: Empirical risk estimates in-sample risk.  
B: Thus the out of sample risk is also small.

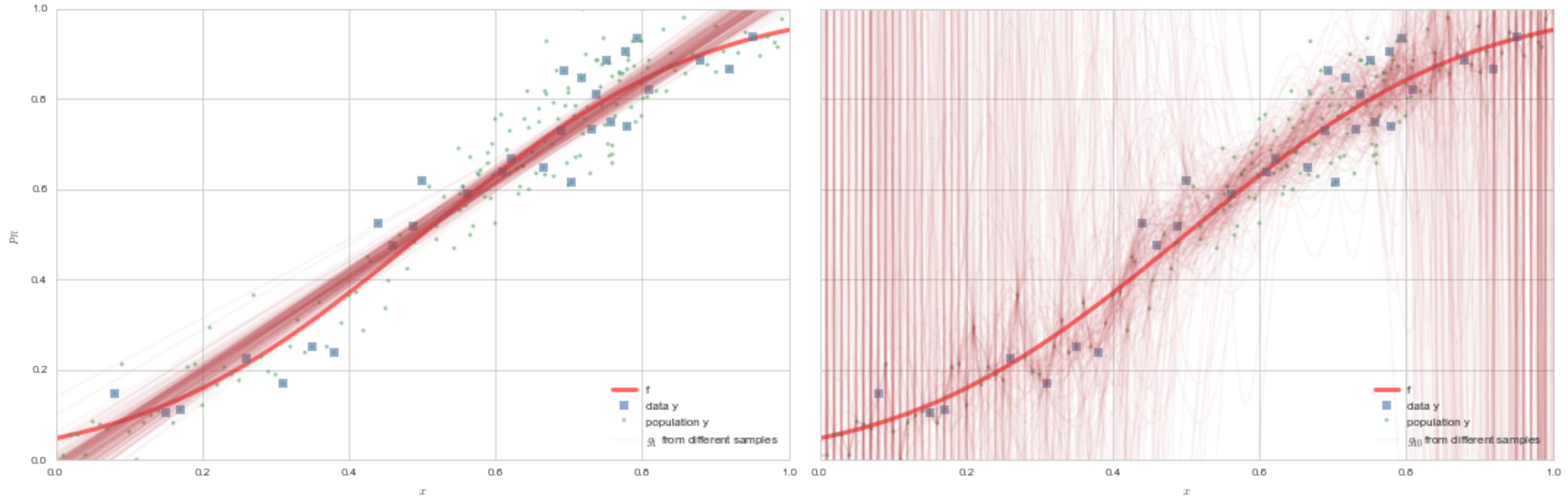
# Bias or Underfitting

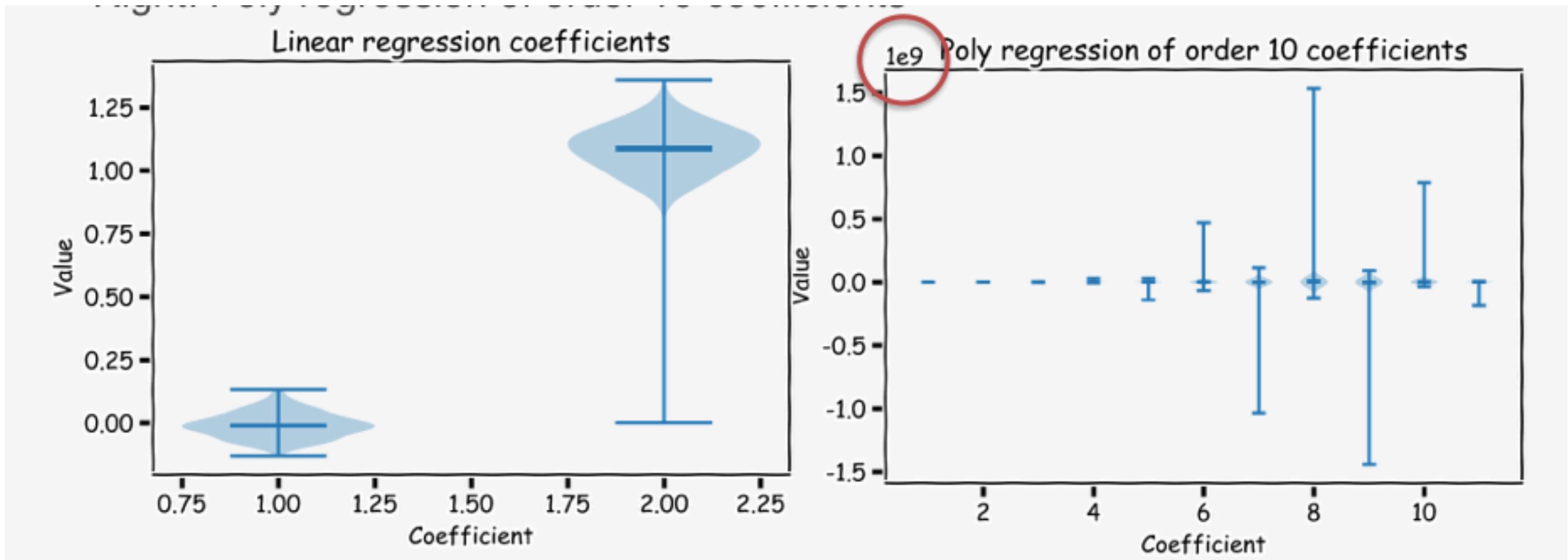


# Data size matters

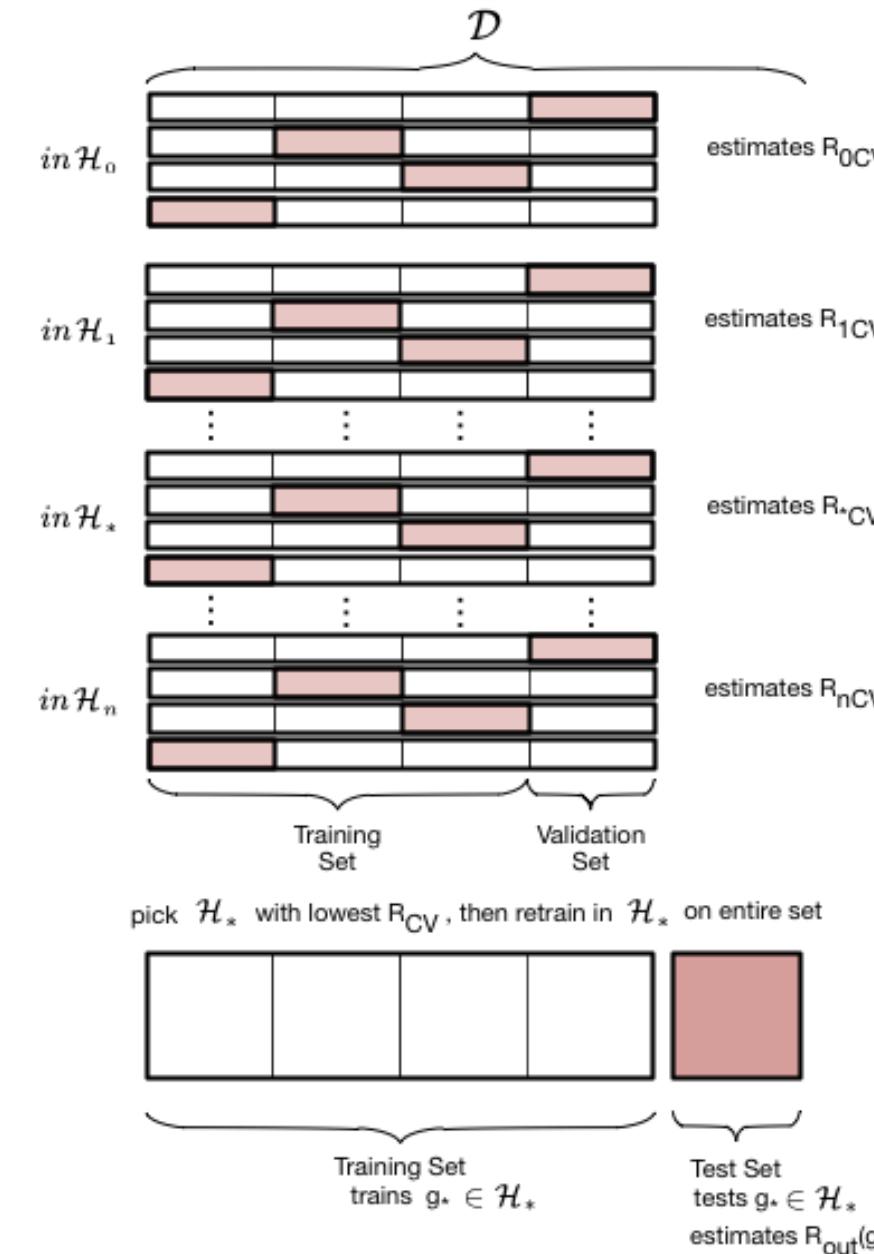
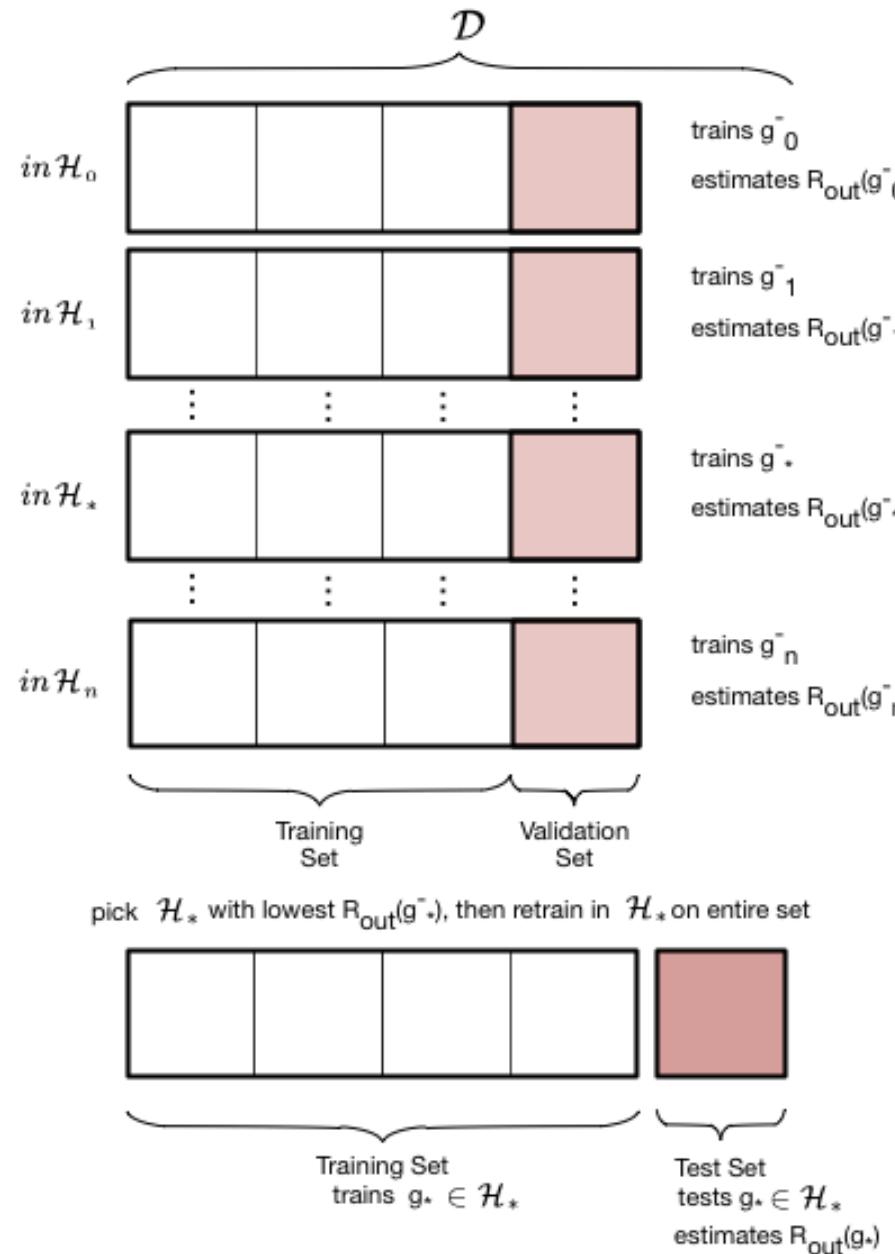


# UNDERFIT(Bias) vs OVERFIT (Variance)

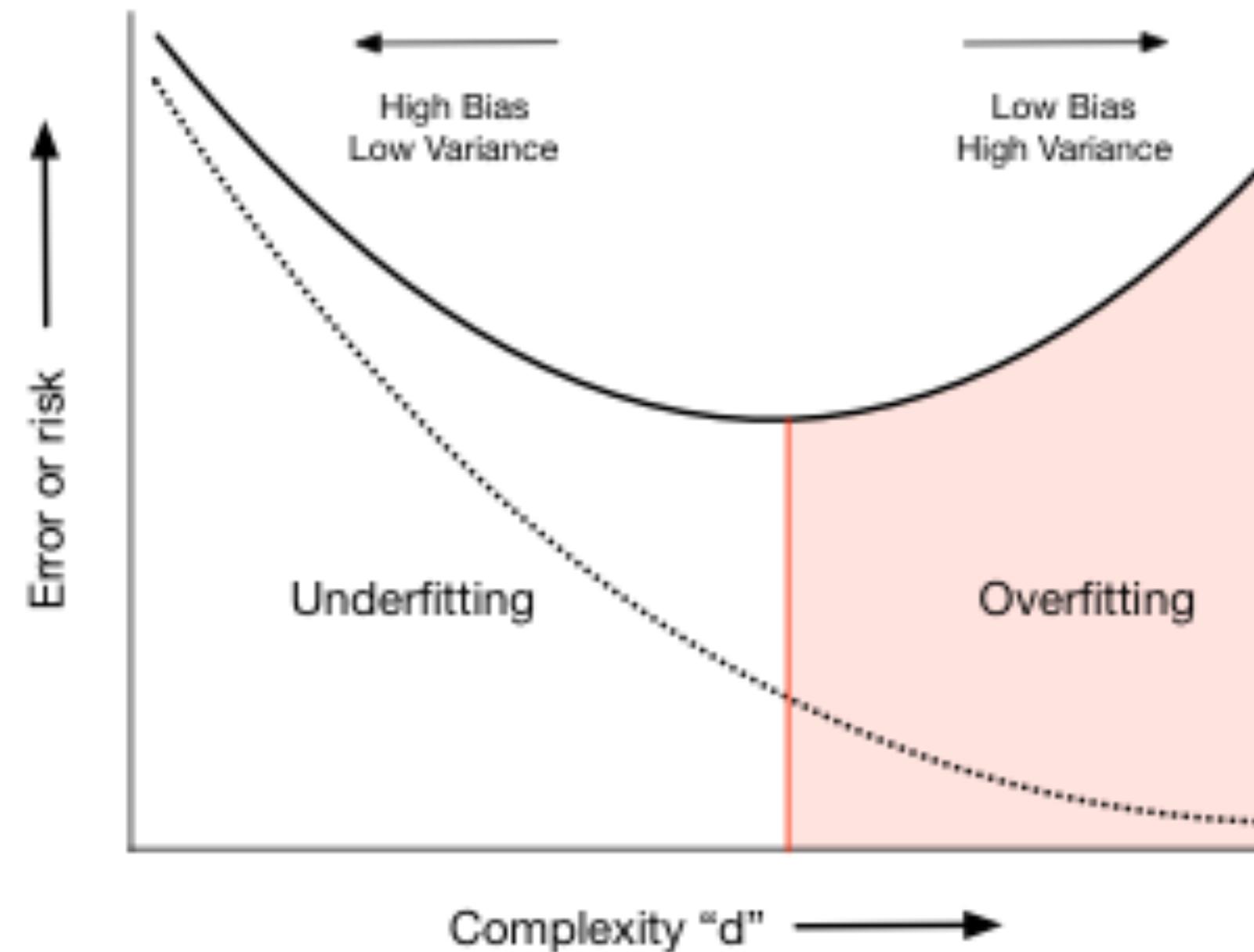




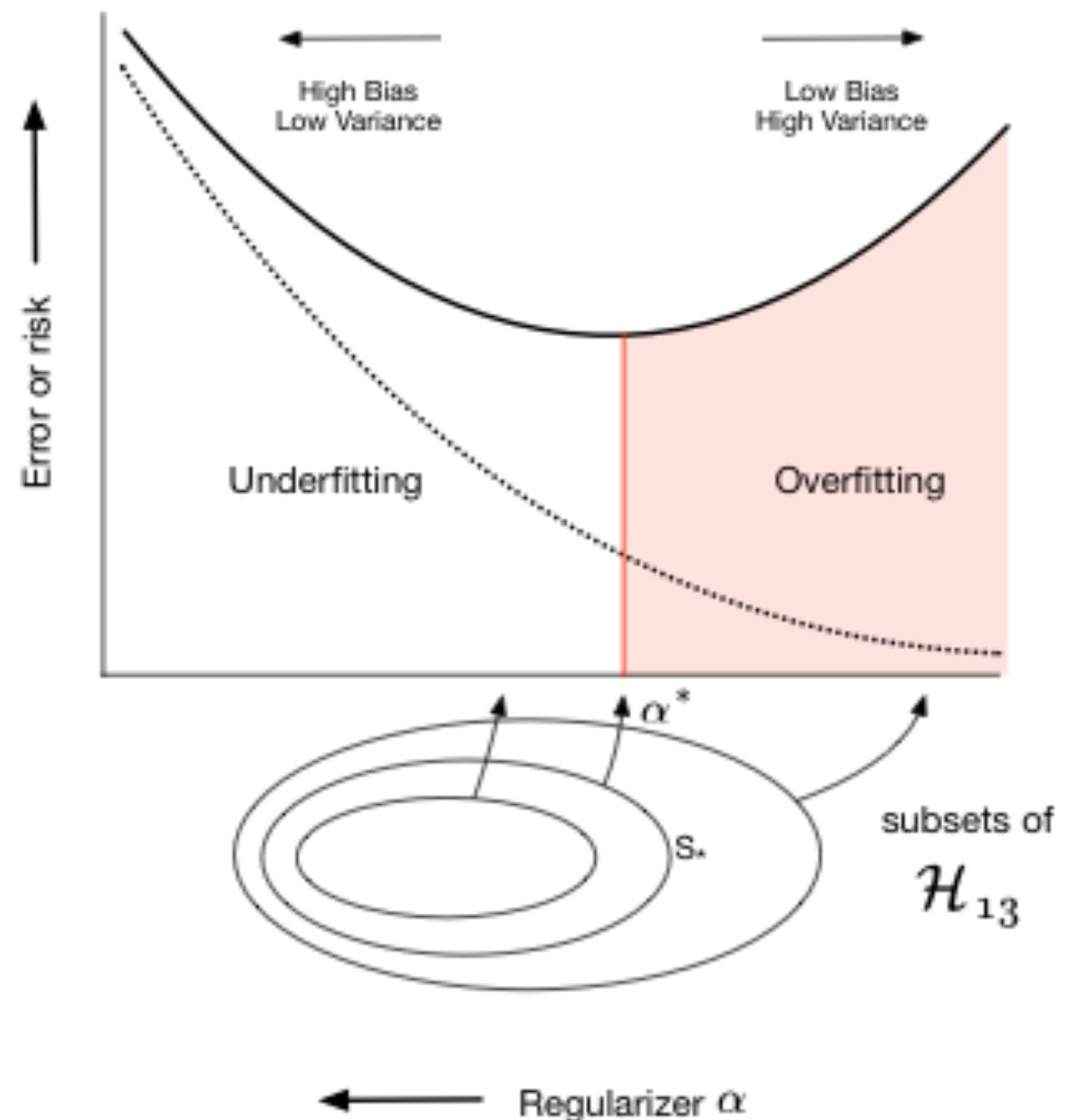
# Validation, X-Val



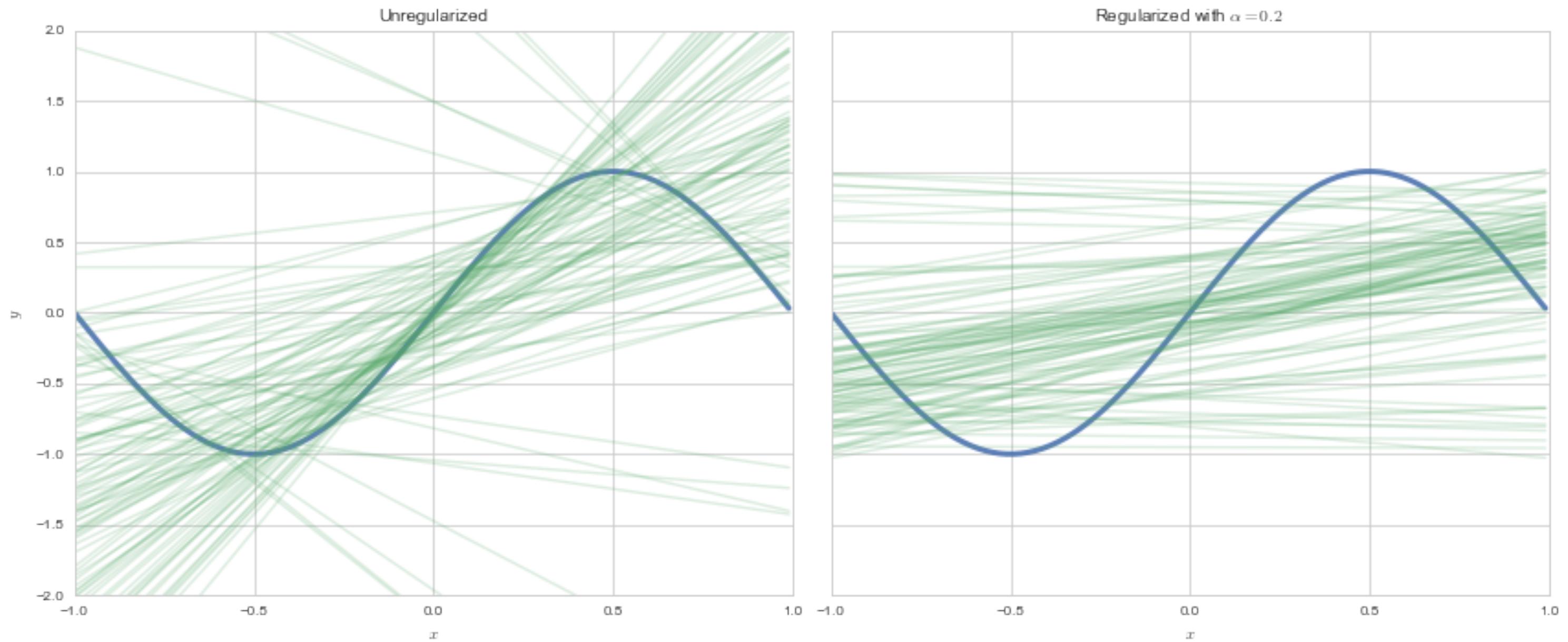
# Dont Overfit



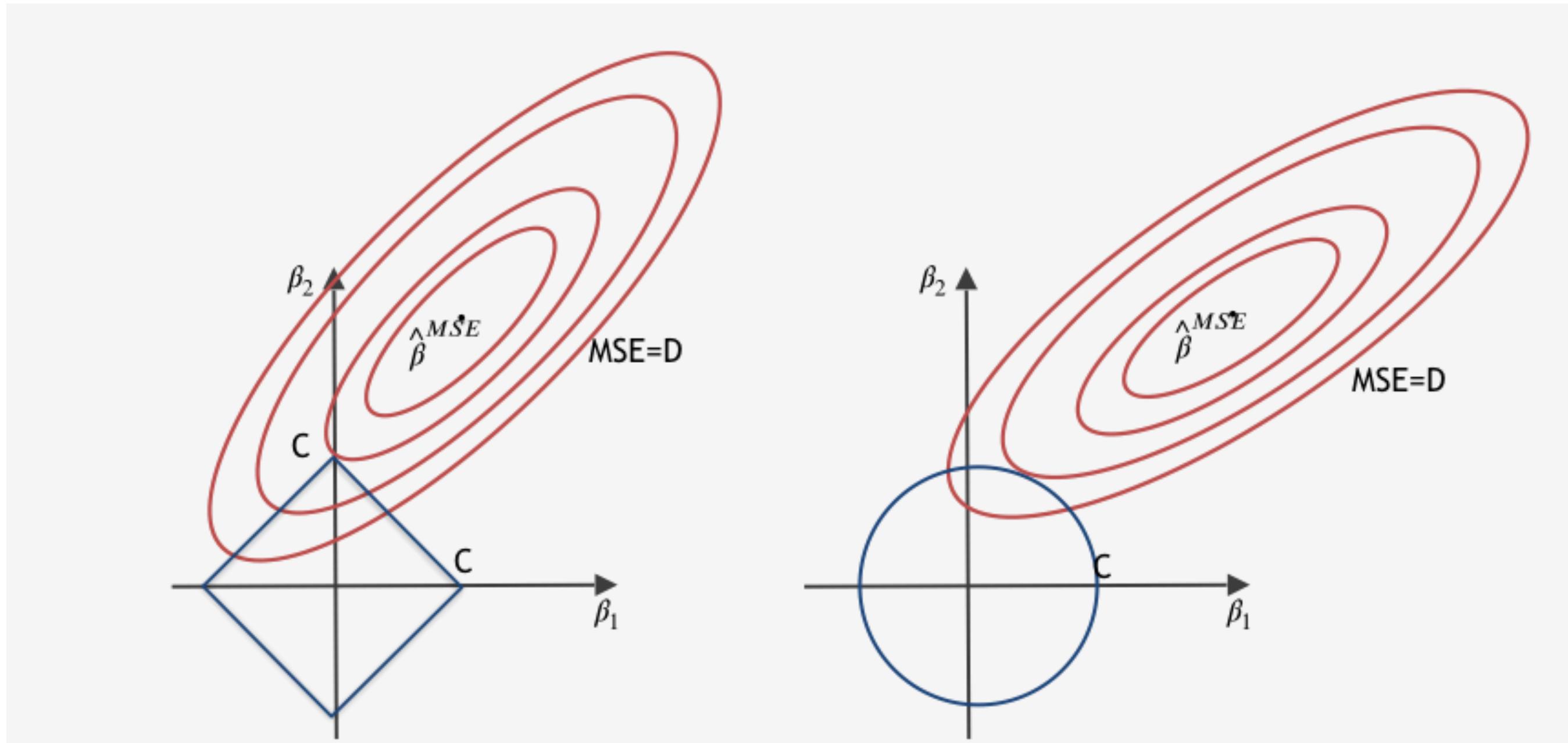
# Regularization



# Make em behave!

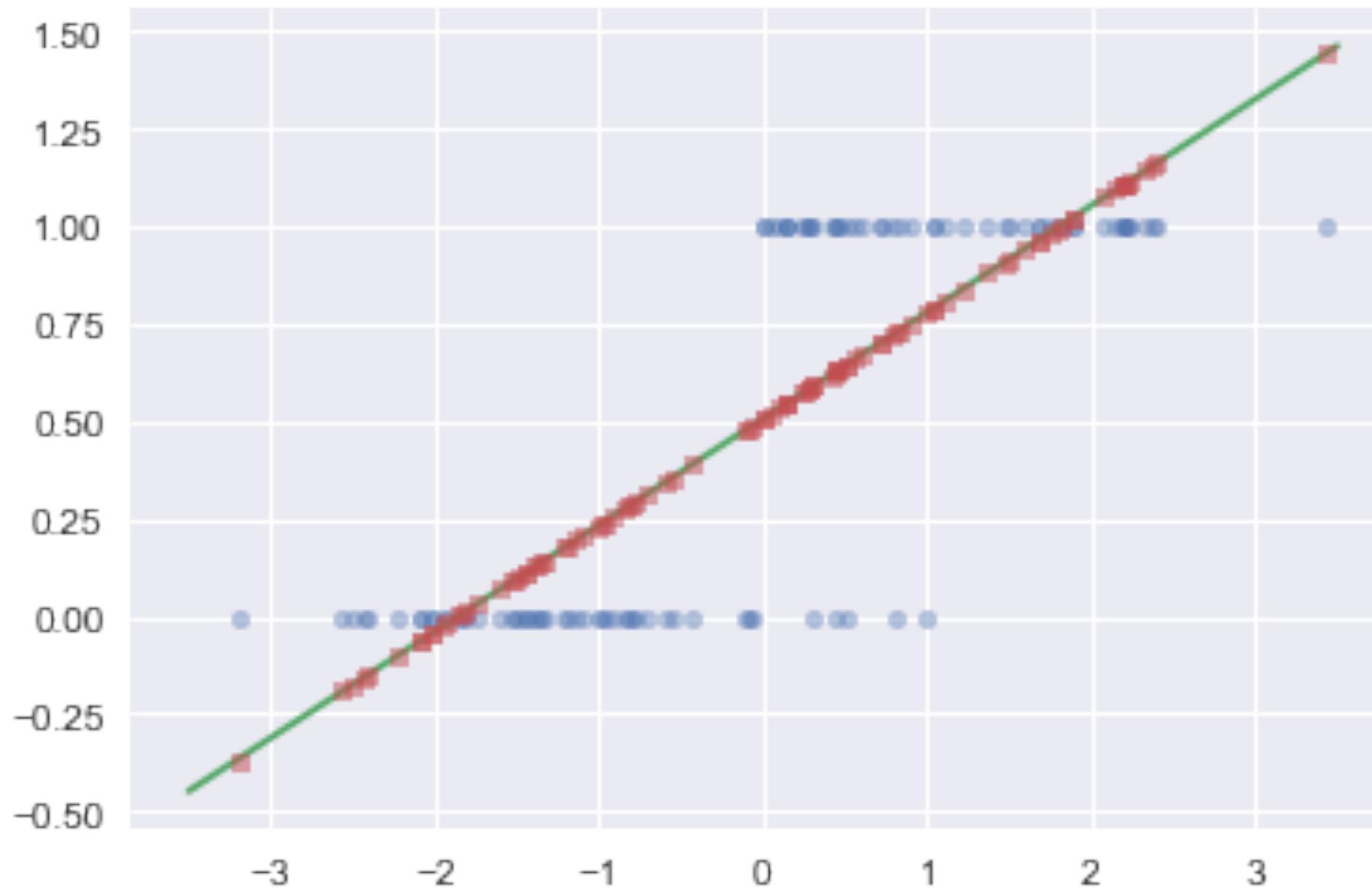


# Geometry of regularization

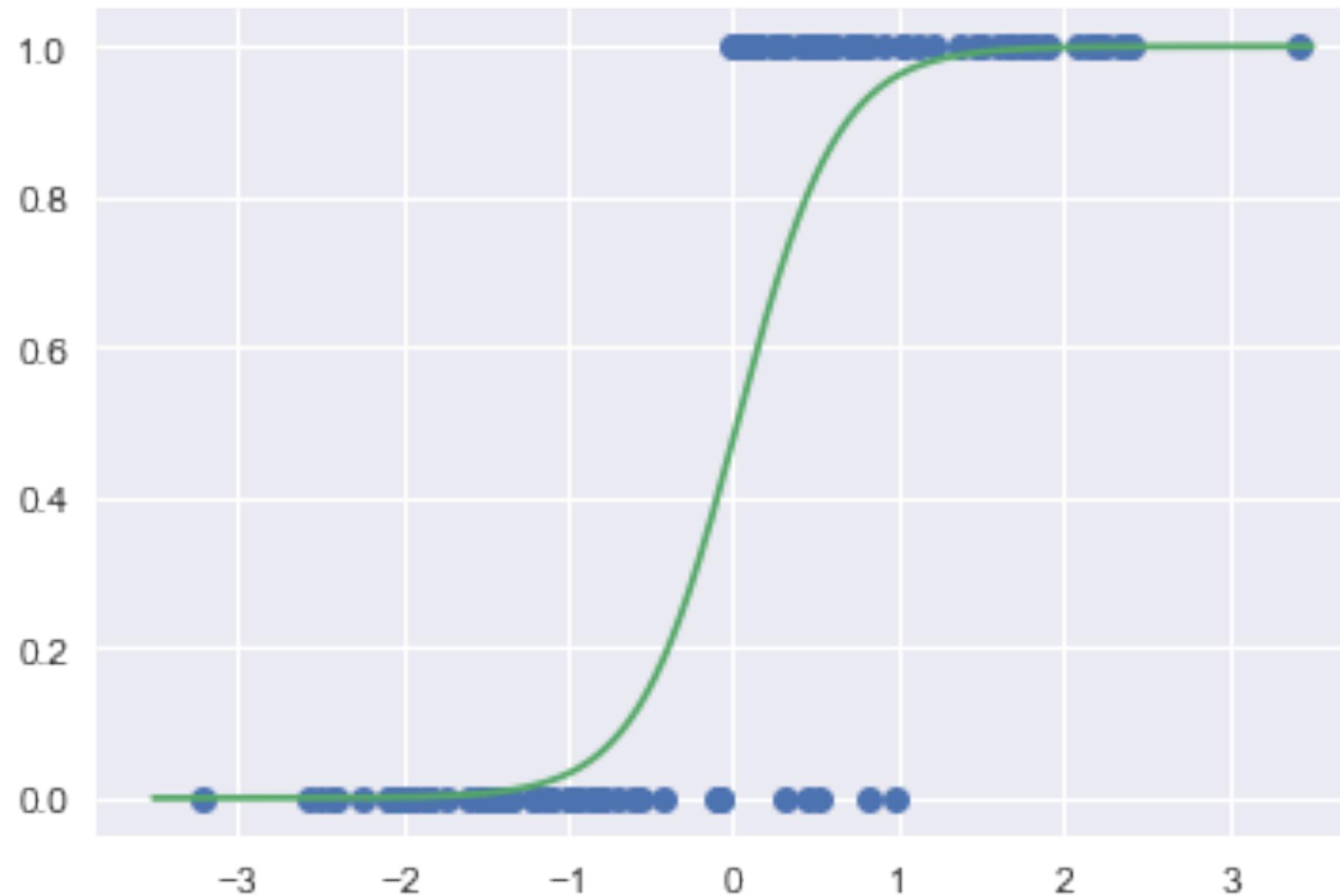


# Classification

# With Linear Regression



# With Logistic Regression

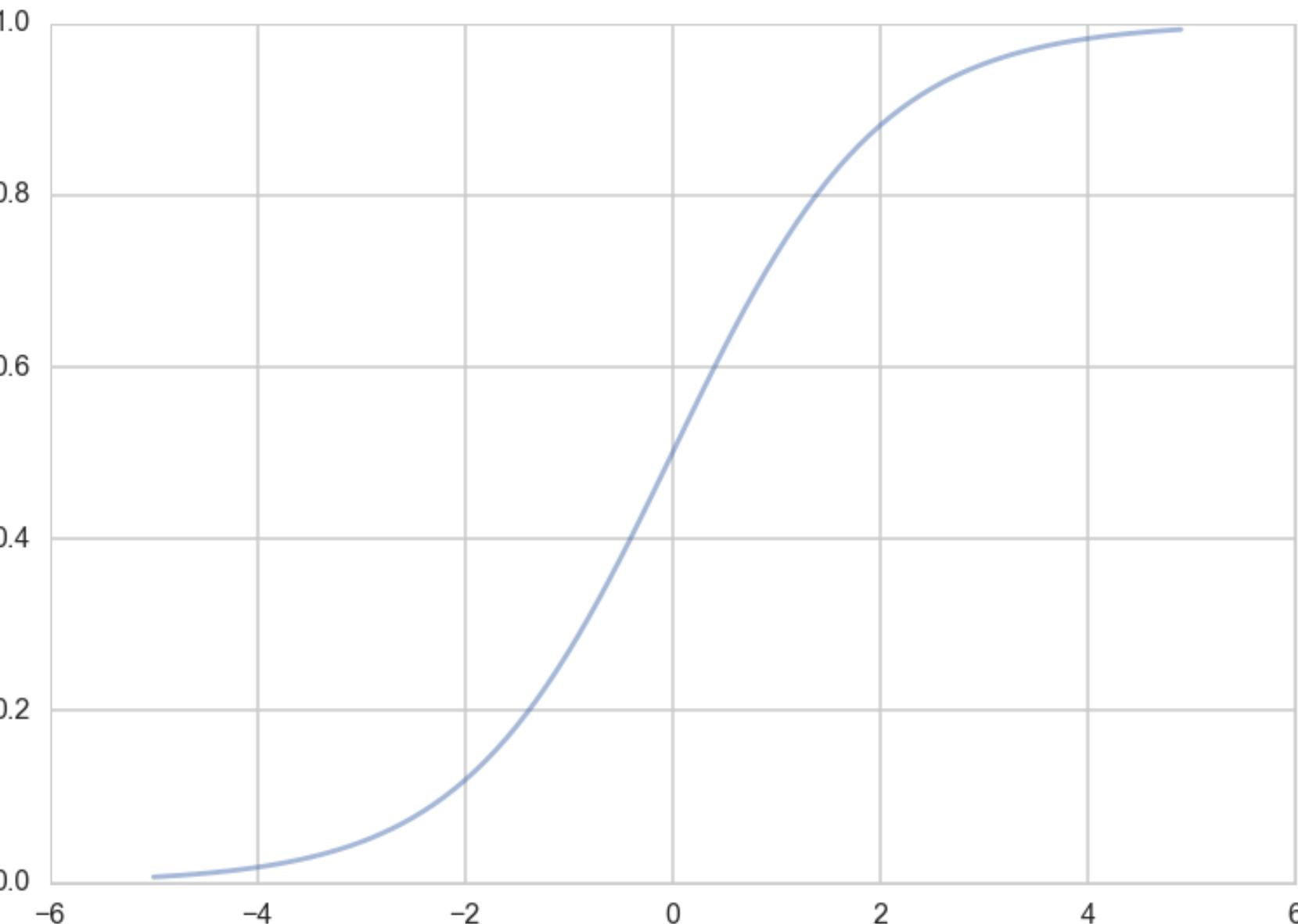


# Sigmoid function

This function is plotted below:

```
h = lambda z: 1./(1+np.exp(-z))  
zs=np.arange(-5,5,0.1)  
plt.plot(zs, h(zs), alpha=0.5);
```

Identify:  $z = \mathbf{w} \cdot \mathbf{x}$  and  $h(\mathbf{w} \cdot \mathbf{x})$  with the probability that the sample is a '1' ( $y = 1$ ).



Then, the conditional probabilities of  $y = 1$  or  $y = 0$  given a particular sample's features  $\mathbf{x}$  are:

$$P(y = 1|\mathbf{x}) = h(\mathbf{w} \cdot \mathbf{x})$$

$$P(y = 0|\mathbf{x}) = 1 - h(\mathbf{w} \cdot \mathbf{x}).$$

These two can be written together as

$$P(y|\mathbf{x}, \mathbf{w}) = h(\mathbf{w} \cdot \mathbf{x})^y (1 - h(\mathbf{w} \cdot \mathbf{x}))^{(1-y)}$$

BERNOULLI!!

Multiplying over the samples we get:

$$P(y|\mathbf{x}, \mathbf{w}) = P(\{y_i\}|\{\mathbf{x}_i\}, \mathbf{w}) = \prod_{y_i \in \mathcal{D}} P(y_i|\mathbf{x}_i, \mathbf{w}) = \prod_{y_i \in \mathcal{D}} h(\mathbf{w} \cdot \mathbf{x}_i)^{y_i} (1 - h(\mathbf{w} \cdot \mathbf{x}_i))^{(1-y_i)}$$

Indeed its important to realize that a particular sample can be thought of as a draw from some "true" probability distribution.

**maximum likelihood** estimation maximises the **likelihood of the sample  $y$** , or alternately the log-likelihood,

$$\mathcal{L} = P(y | \mathbf{x}, \mathbf{w}). \text{ OR } \ell = \log(P(y | \mathbf{x}, \mathbf{w}))$$

Thus

$$\begin{aligned}\ell &= \log \left( \prod_{y_i \in \mathcal{D}} h(\mathbf{w} \cdot \mathbf{x}_i)^{y_i} (1 - h(\mathbf{w} \cdot \mathbf{x}_i))^{(1-y_i)} \right) \\ &= \sum_{y_i \in \mathcal{D}} \log \left( h(\mathbf{w} \cdot \mathbf{x}_i)^{y_i} (1 - h(\mathbf{w} \cdot \mathbf{x}_i))^{(1-y_i)} \right) \\ &= \sum_{y_i \in \mathcal{D}} \log h(\mathbf{w} \cdot \mathbf{x}_i)^{y_i} + \log (1 - h(\mathbf{w} \cdot \mathbf{x}_i))^{(1-y_i)} \\ &= \sum_{y_i \in \mathcal{D}} (y_i \log(h(\mathbf{w} \cdot \mathbf{x})) + (1 - y_i) \log(1 - h(\mathbf{w} \cdot \mathbf{x})))\end{aligned}$$

## Logistic Regression: NLL

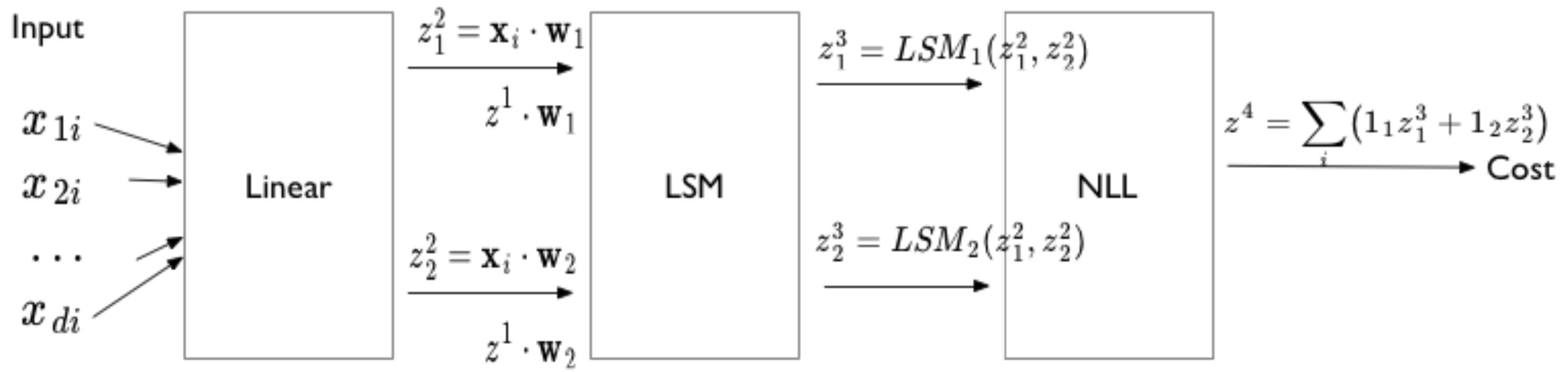
The negative of this log likelihood (NLL), also called *cross-entropy*.

$$NLL = - \sum_{y_i \in \mathcal{D}} (y_i \log(h(\mathbf{w} \cdot \mathbf{x})) + (1 - y_i) \log(1 - h(\mathbf{w} \cdot \mathbf{x})))$$

Gradient:  $\nabla_{\mathbf{w}} NLL = \sum_i \mathbf{x}_i^T (p_i - y_i) = \mathbf{X}^T \cdot (\mathbf{p} - \mathbf{w})$

Hessian:  $H = \mathbf{X}^T diag(p_i(1 - p_i)) \mathbf{X}$  positive definite  $\implies$  convex

# Softmax Formulation of Logistic Regression



$$z^1 = \mathbf{x}_i$$

```
(X_train, y_train), (X_test, y_test) = mnist.load_data()

# one hot encode outputs
y_train = to_categorical(y_train)
y_test = to_categorical(y_test)
# create model
model = Sequential()
model.add(Flatten(input_shape=(img_width, img_height)))
model.add(Dense(num_classes, activation='softmax'))
model.compile(loss='categorical_crossentropy', optimizer='adam',
              metrics=['accuracy'])
model.summary()
# Fit the model
model.fit(X_train, y_train, epochs=10, validation_data=(X_test, y_test),
           callbacks=[WandbCallback(data_type="image", labels=labels, save_model=False)])
```

```
from sklearn.model_selection import GridSearchCV

pipeline = make_pipeline(TfidfVectorizer(),
                        LogisticRegression())

grid = GridSearchCV(pipeline,
                     param_grid={'logisticregression__C': [.1, 1, 10, 100]}, cv=5)

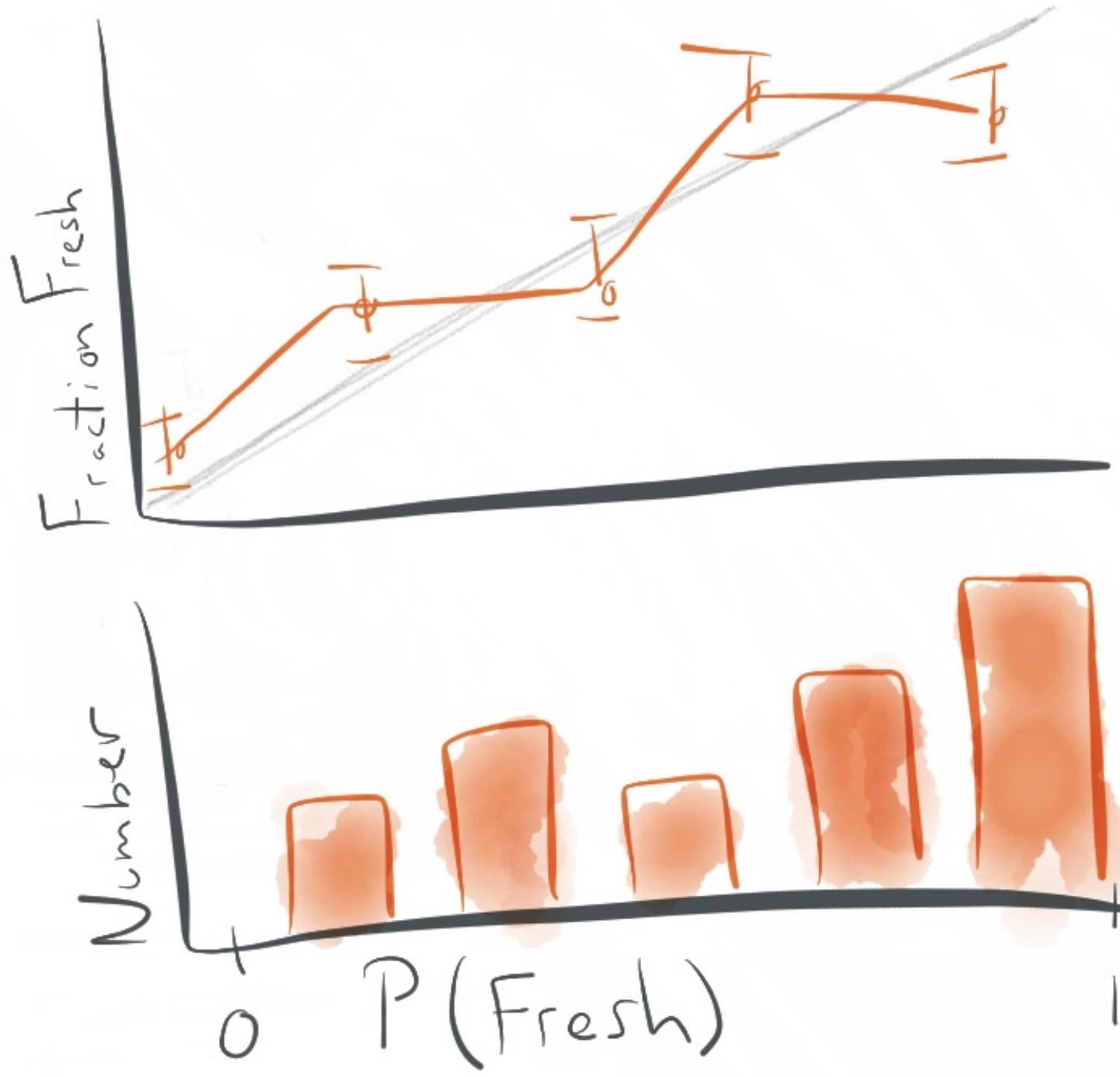
grid.fit(text_train, y_train)
print("Score", grid.score(text_test, y_test))
```

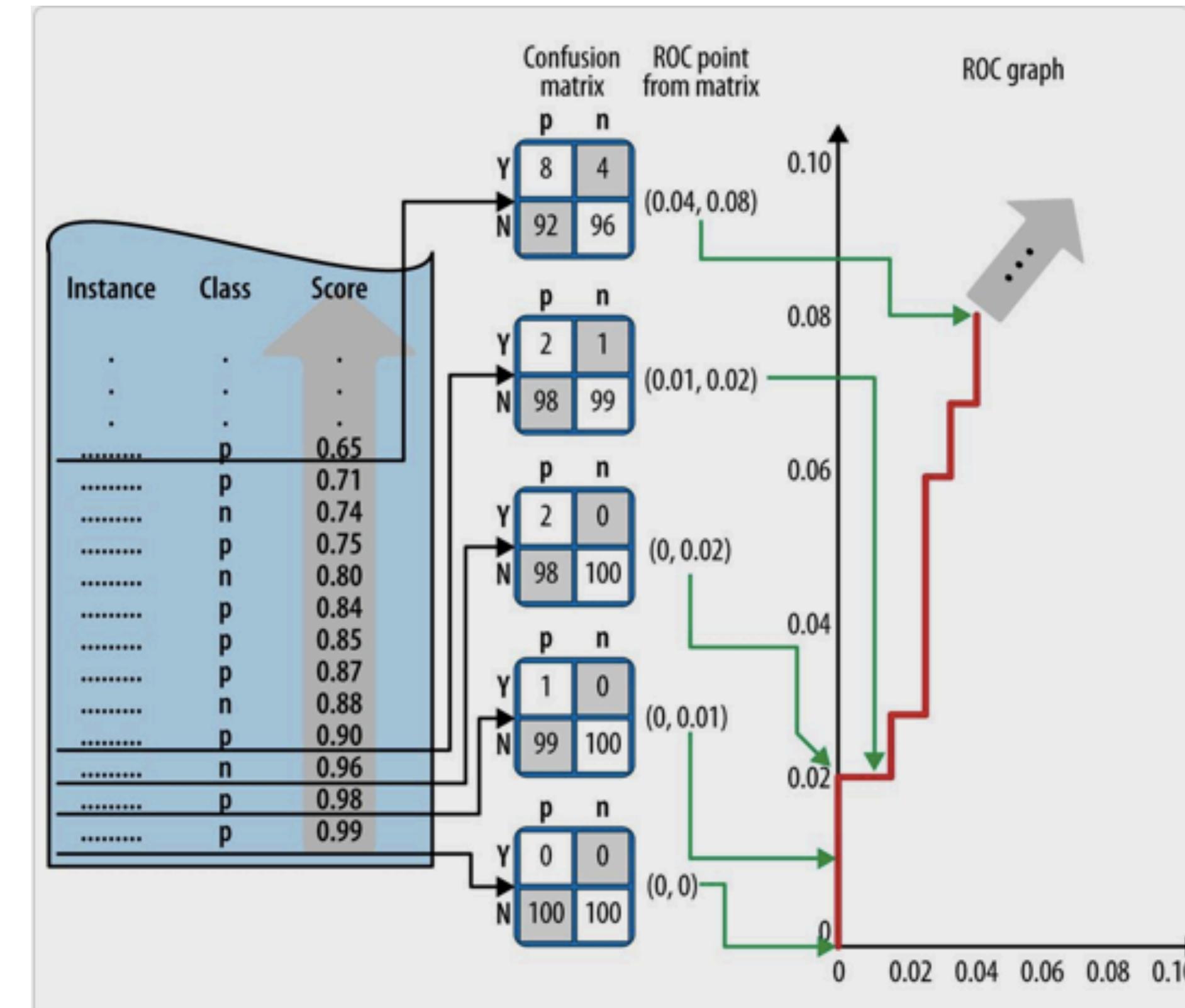
# Metrics and Decision Theory

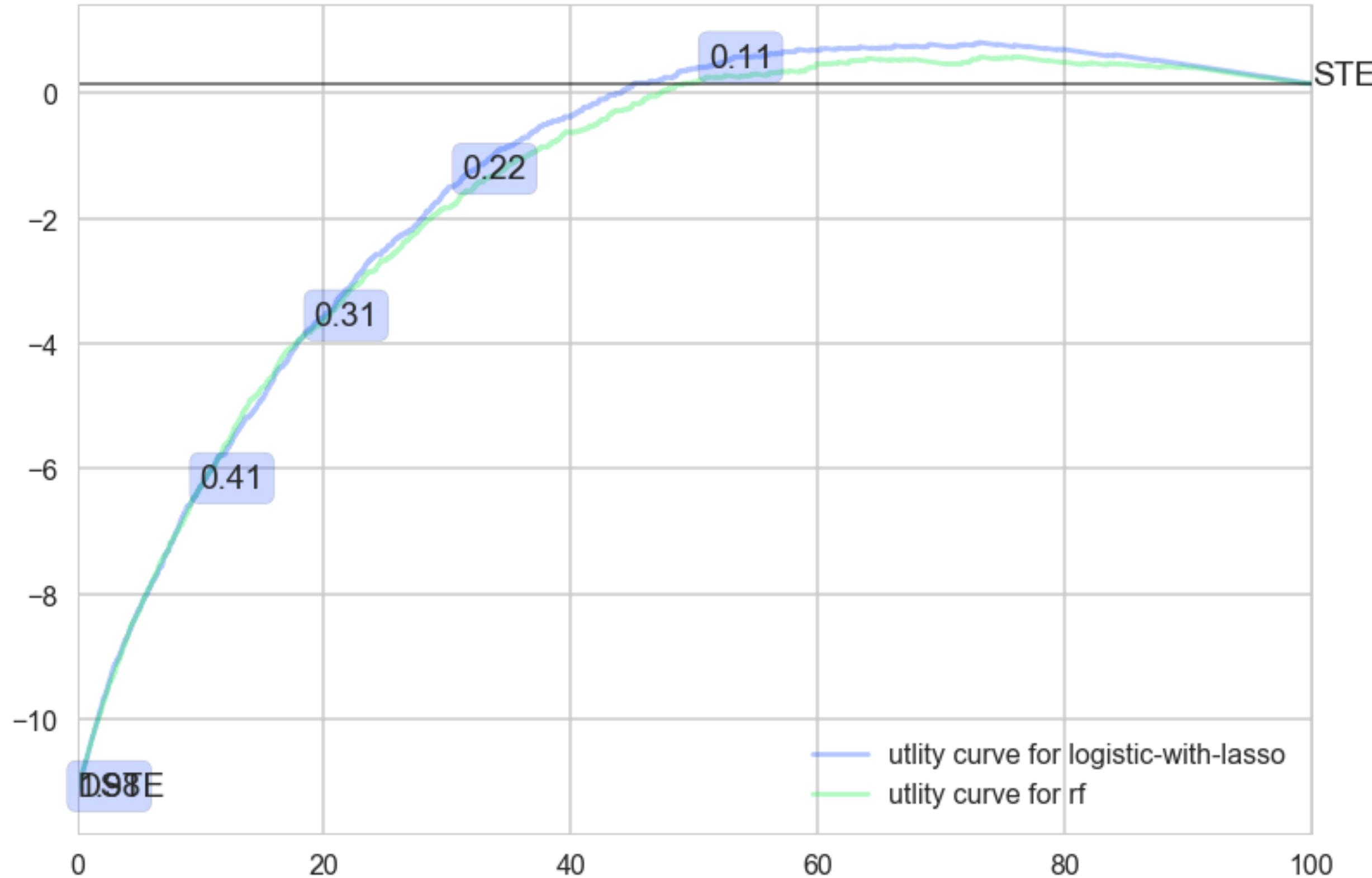
# Confusion Matrix

		Predicted	
		0	1
Observed	0	TN True Negative	FP False Positive
	1	FN False Negative	TP True Positive
		PN Predicted Negative	PP Predicted Positive

		Predicted	
		0	1
Observed	0	TNC True Negative Cost	FPC False Positive Cost
	1	FNC False Negative Cost	TPC True Positive Cost

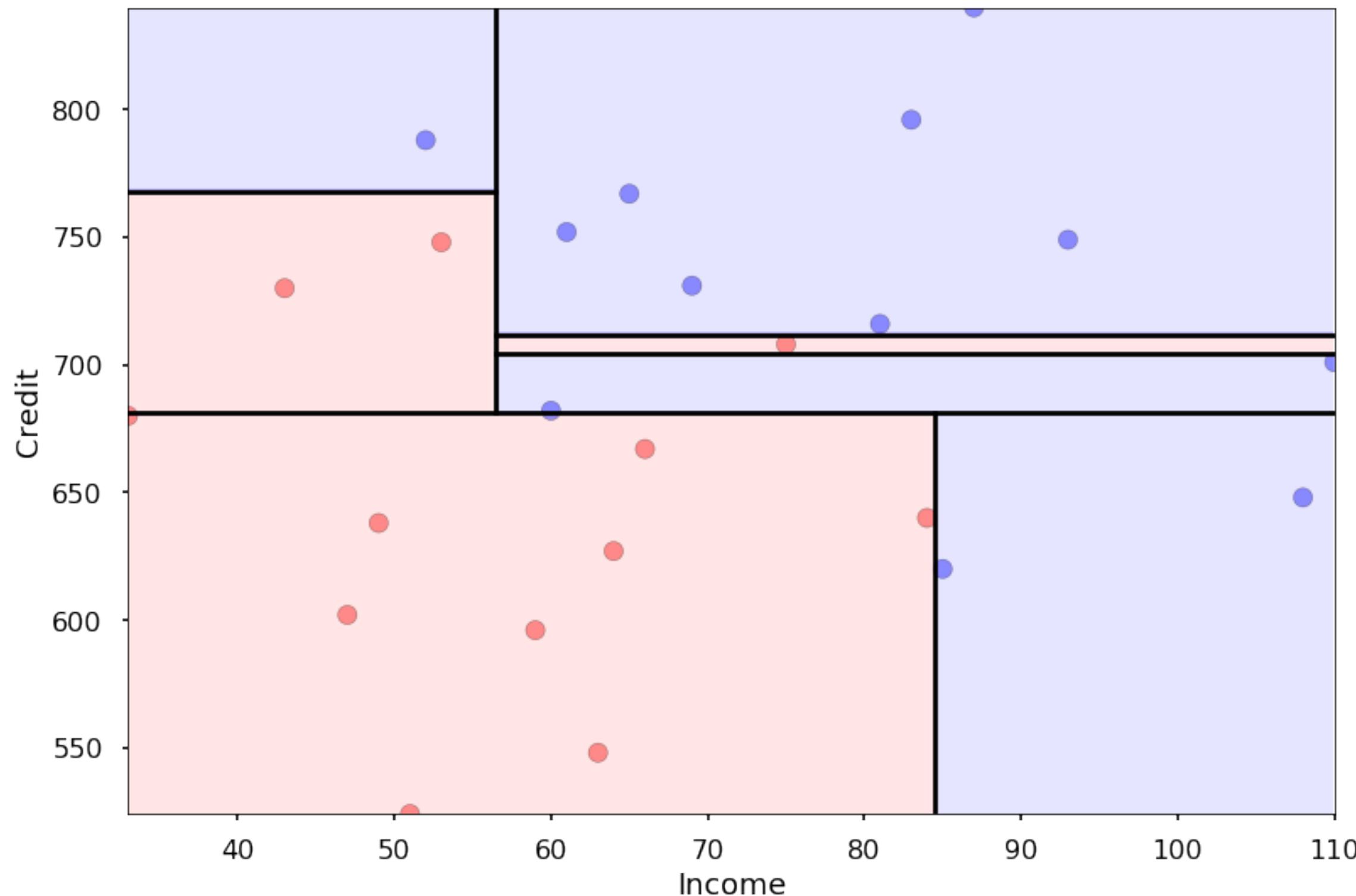


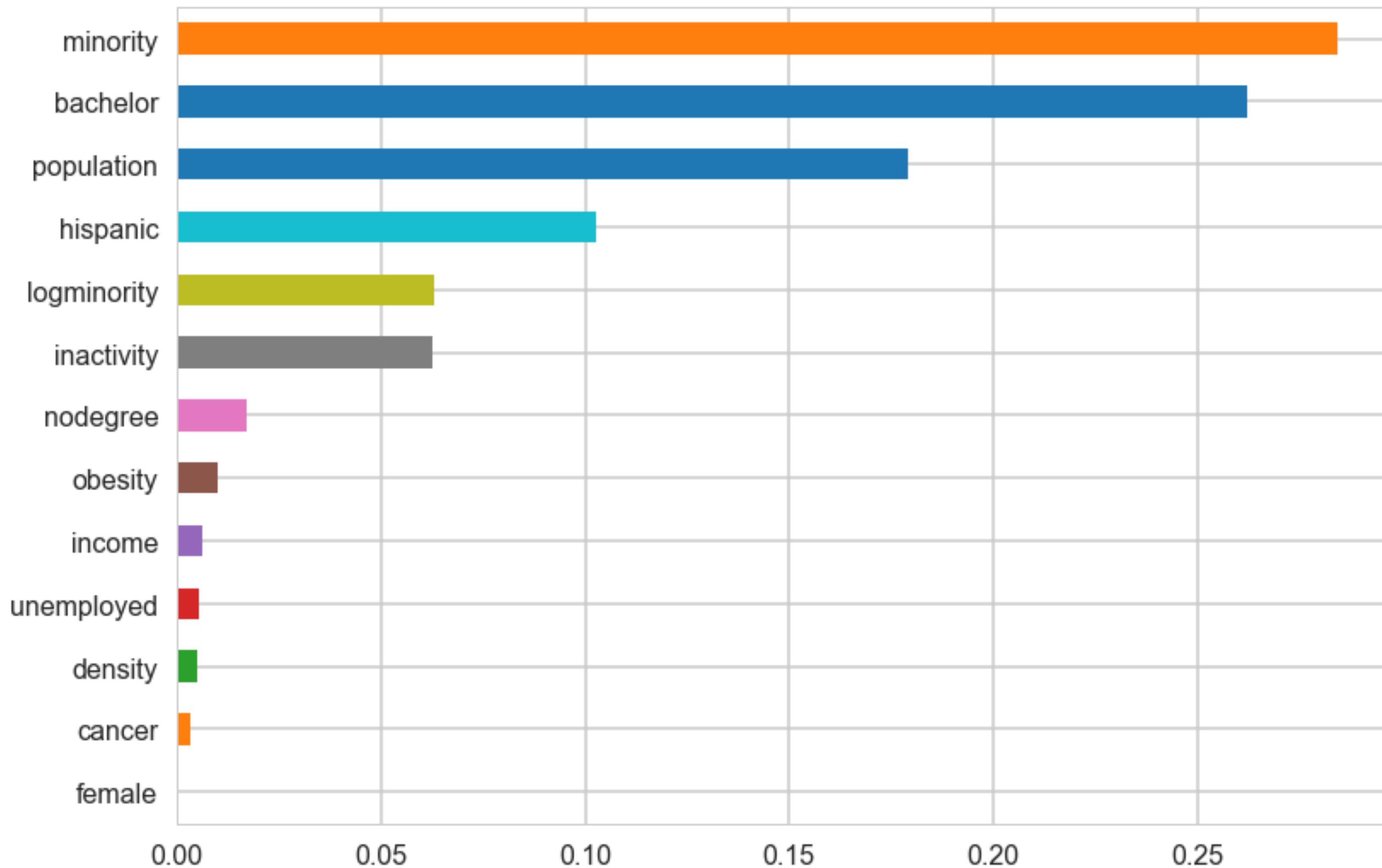


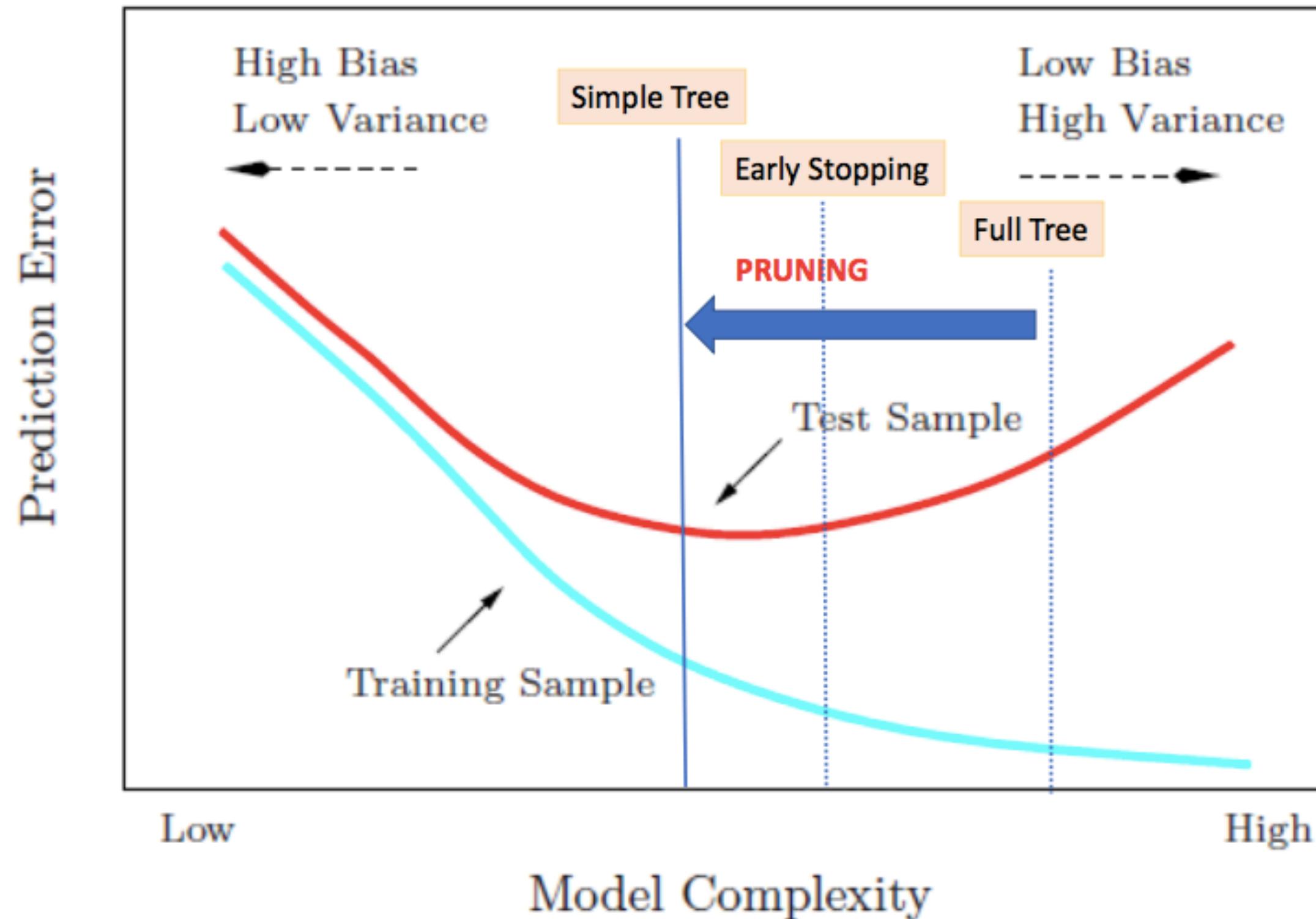


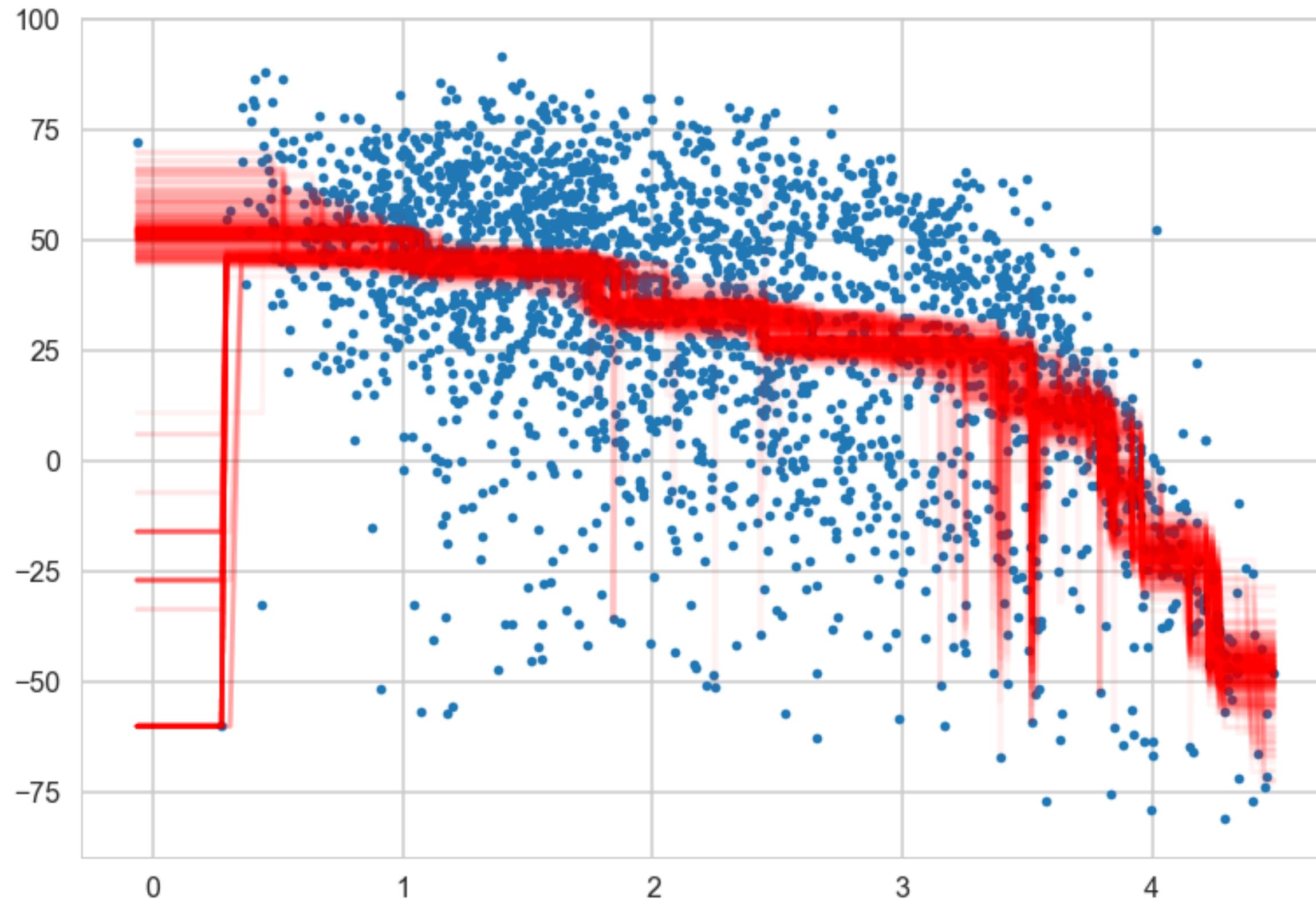
Trees

level 9









# RF

- first do bagging
- then randomly choose features
- Random forest, squeezes out variance
- Big idea is ensembling

# Boosting

- use weak learners
- fit residuals
- gradient descent in "data" space
- inspiration for resnets

neural

Networks

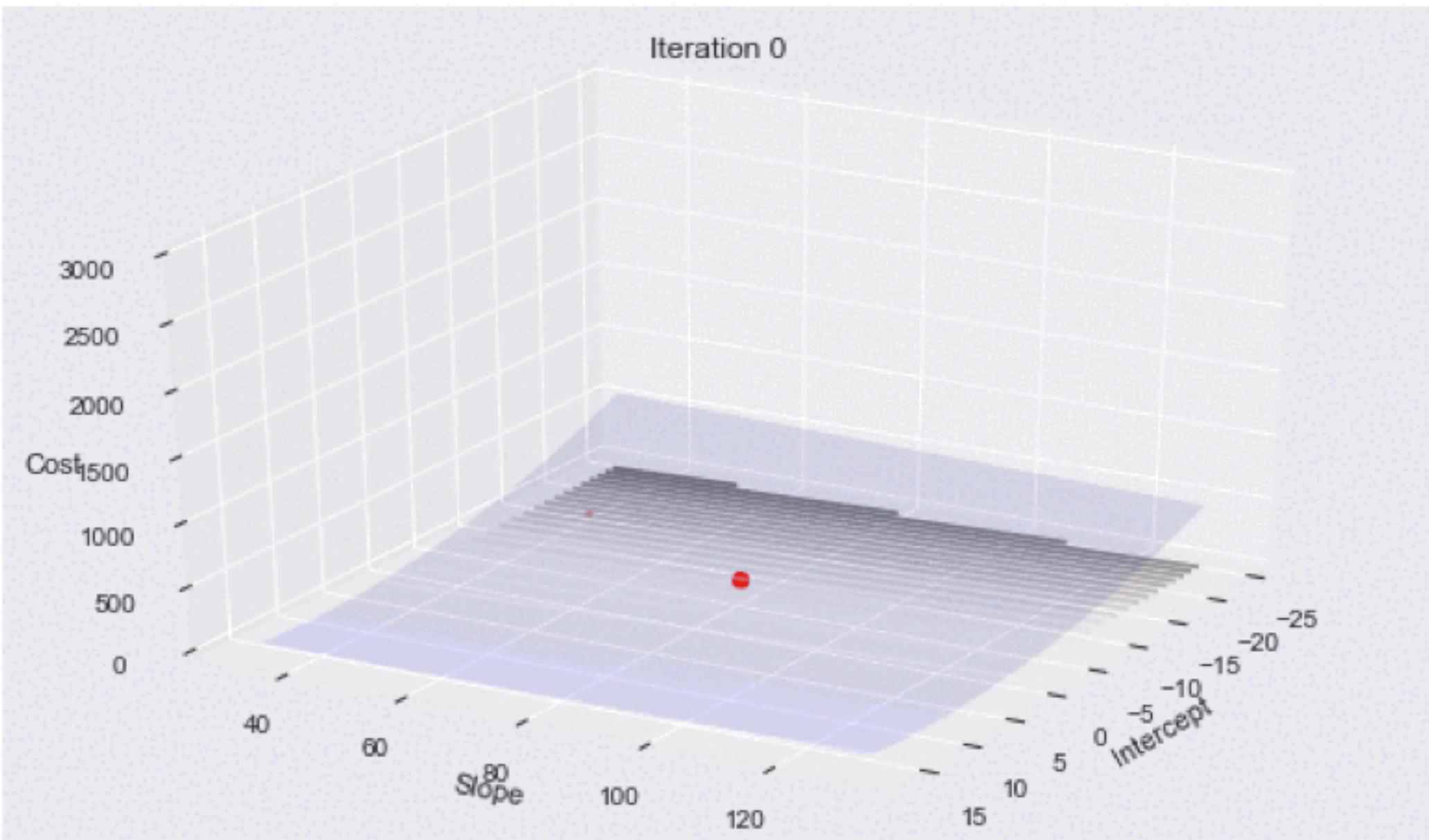
# Stochastic Gradient Descent

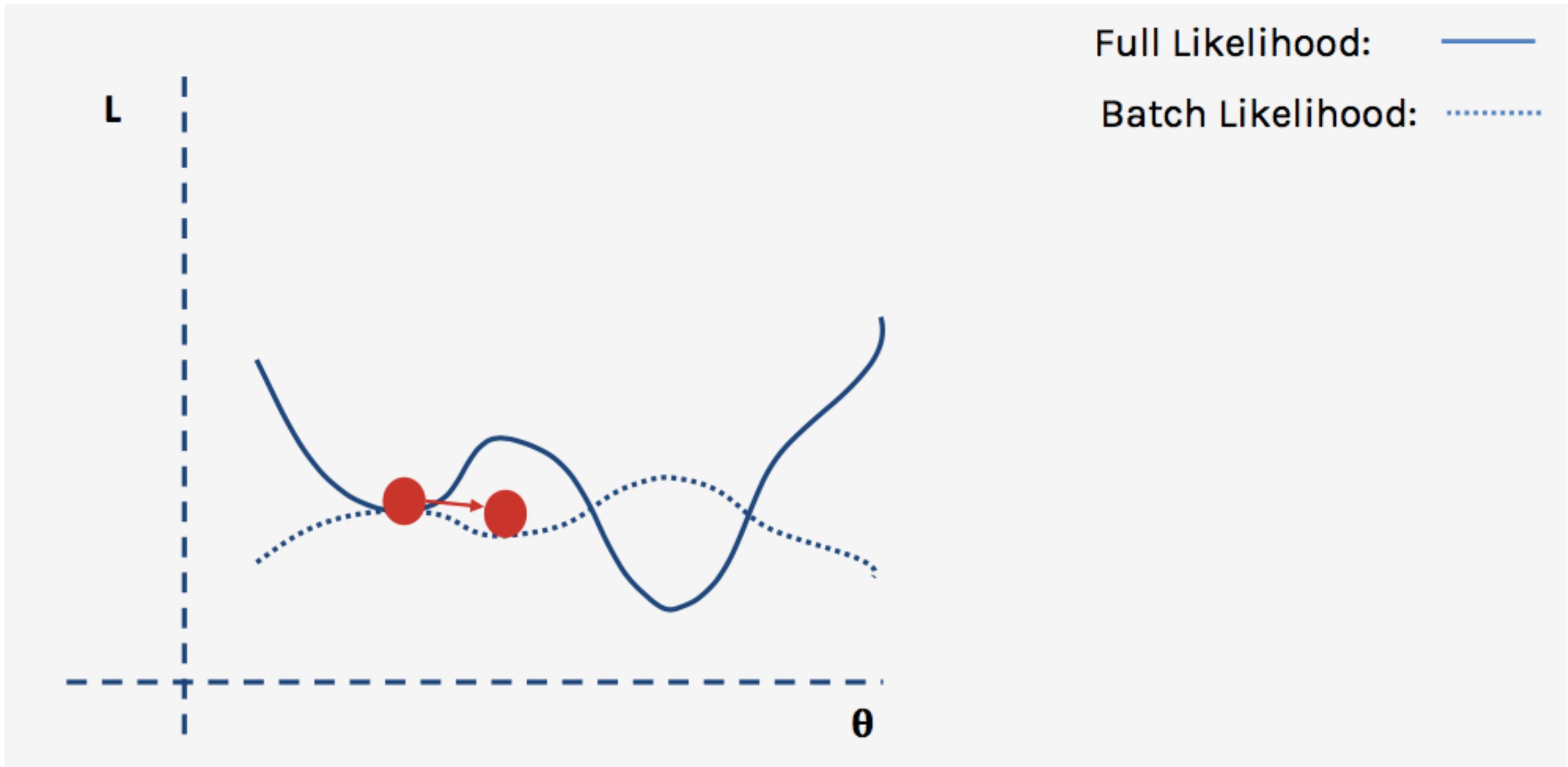
$$\theta := \theta - \alpha \nabla_{\theta} J_i(\theta)$$

ONE POINT AT A TIME

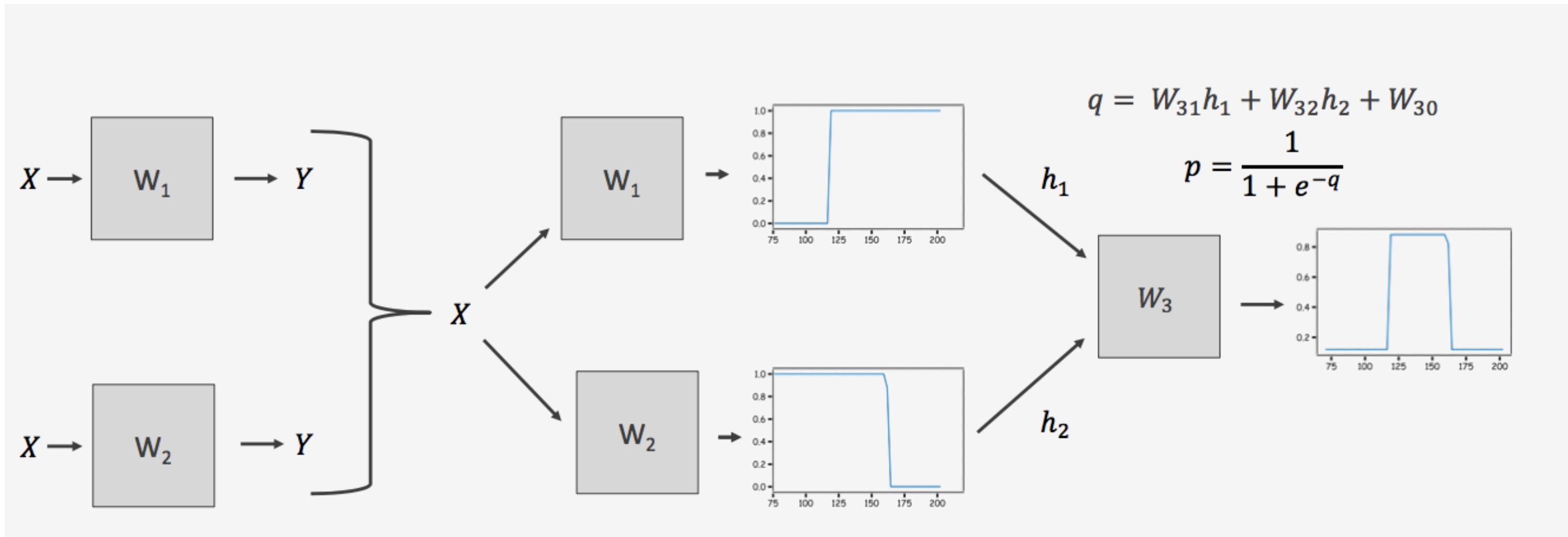
```
for i in range(nb_epochs):
    np.random.shuffle(data)
    for example in data:
        params_grad = evaluate_gradient(loss_function, example, params)
        params = params - learning_rate * params_grad
```

Mini-Batch: do some at a time





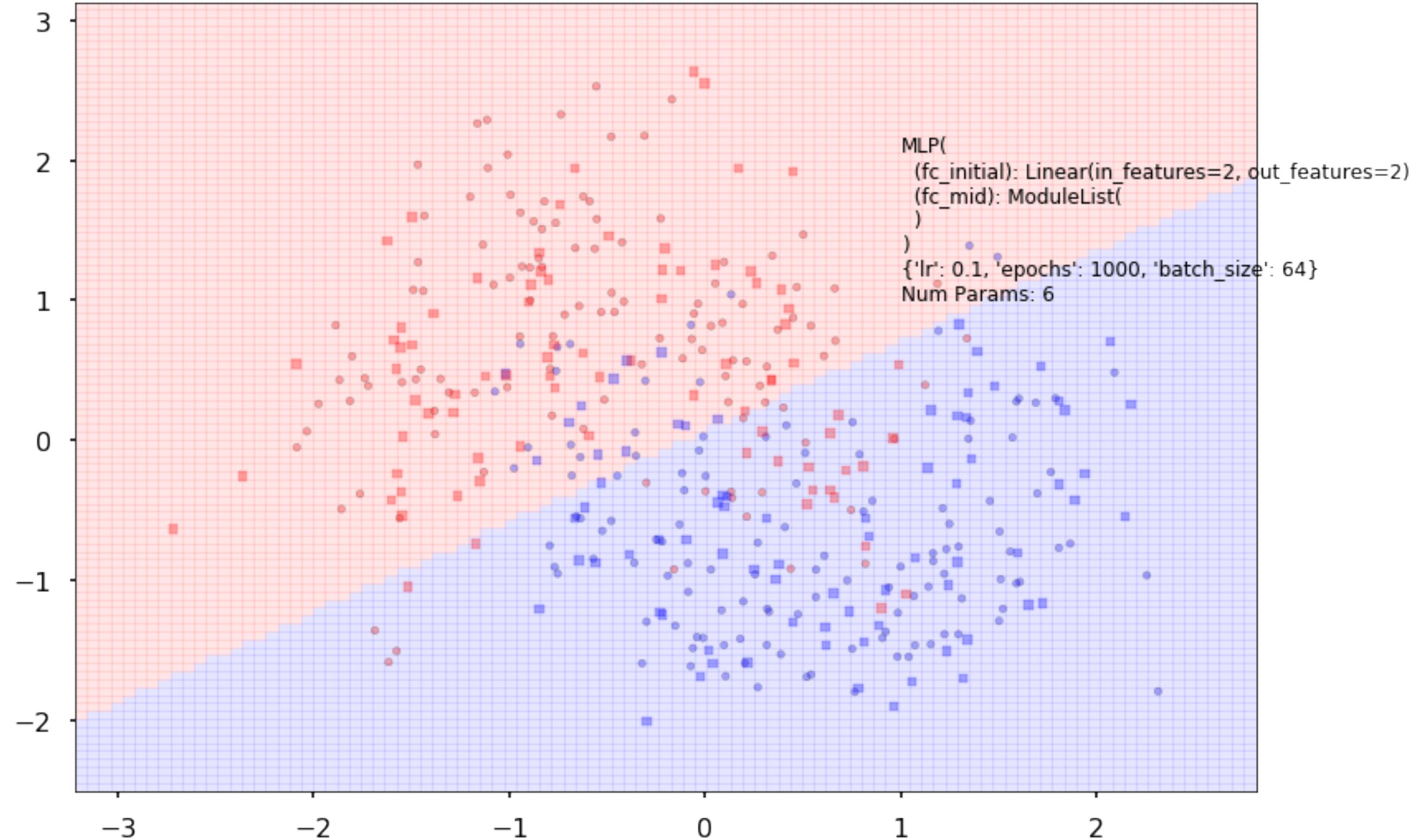
# Basic Idea: Universal approx by combining nonlinears

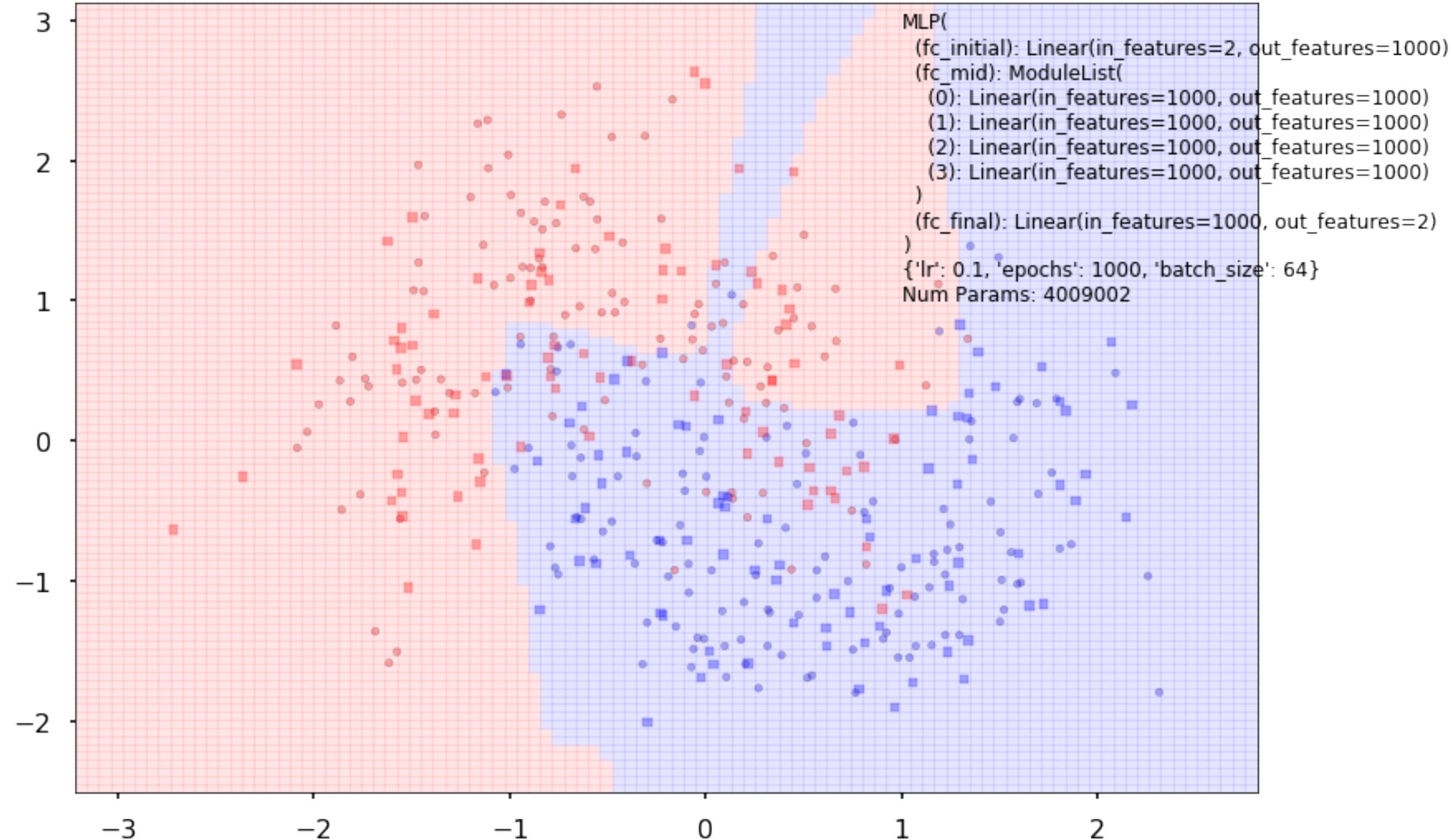


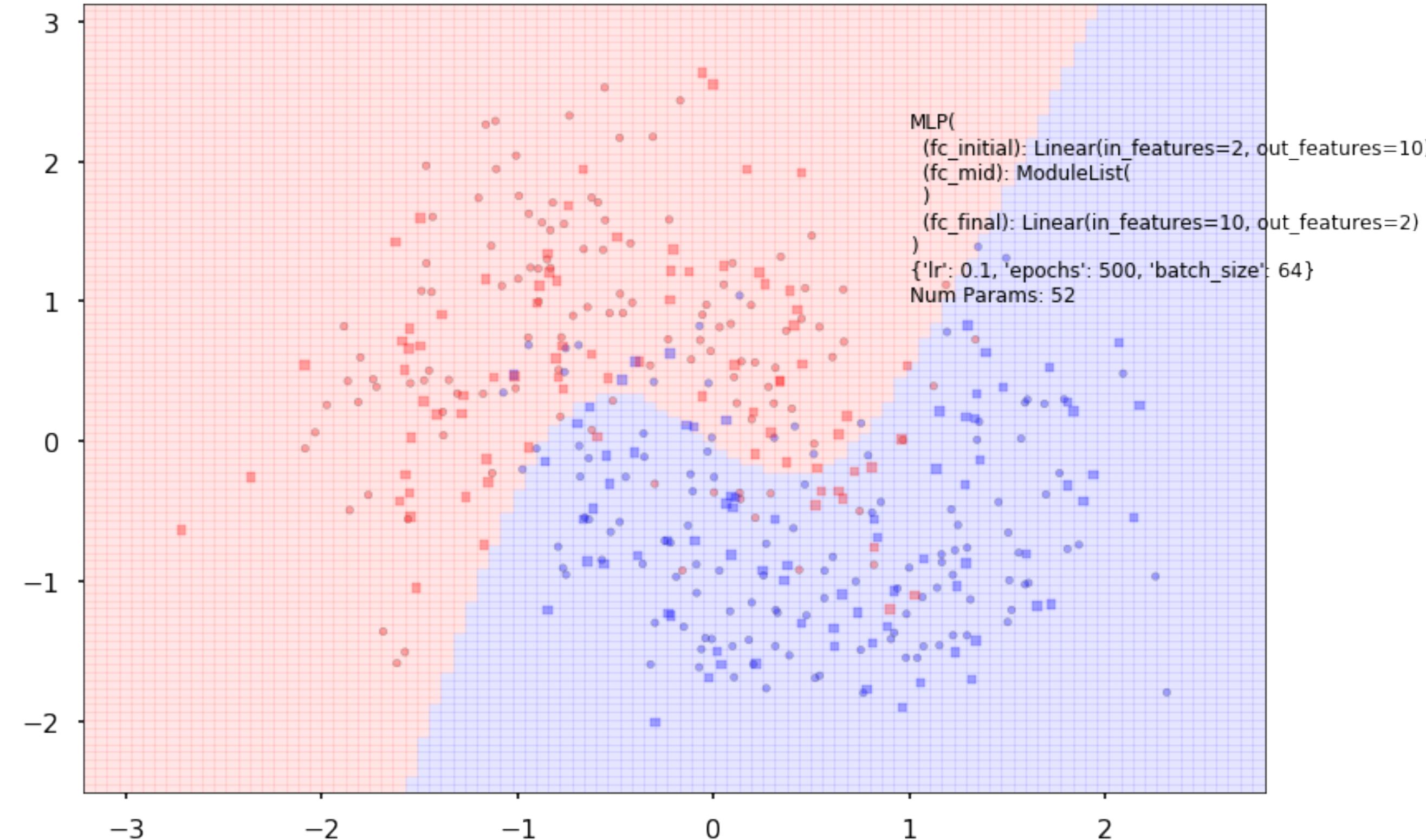
# Multi Layer Perceptrons

```
# create model
model = Sequential()
model.add(Flatten(input_shape=(img_width, img_height)))
model.add(Dense(config.hidden_nodes, activation='relu'))
model.add(Dense(num_classes, activation='softmax'))
model.compile(loss='categorical_crossentropy', optimizer=config.optimizer,
               metrics=['accuracy'])

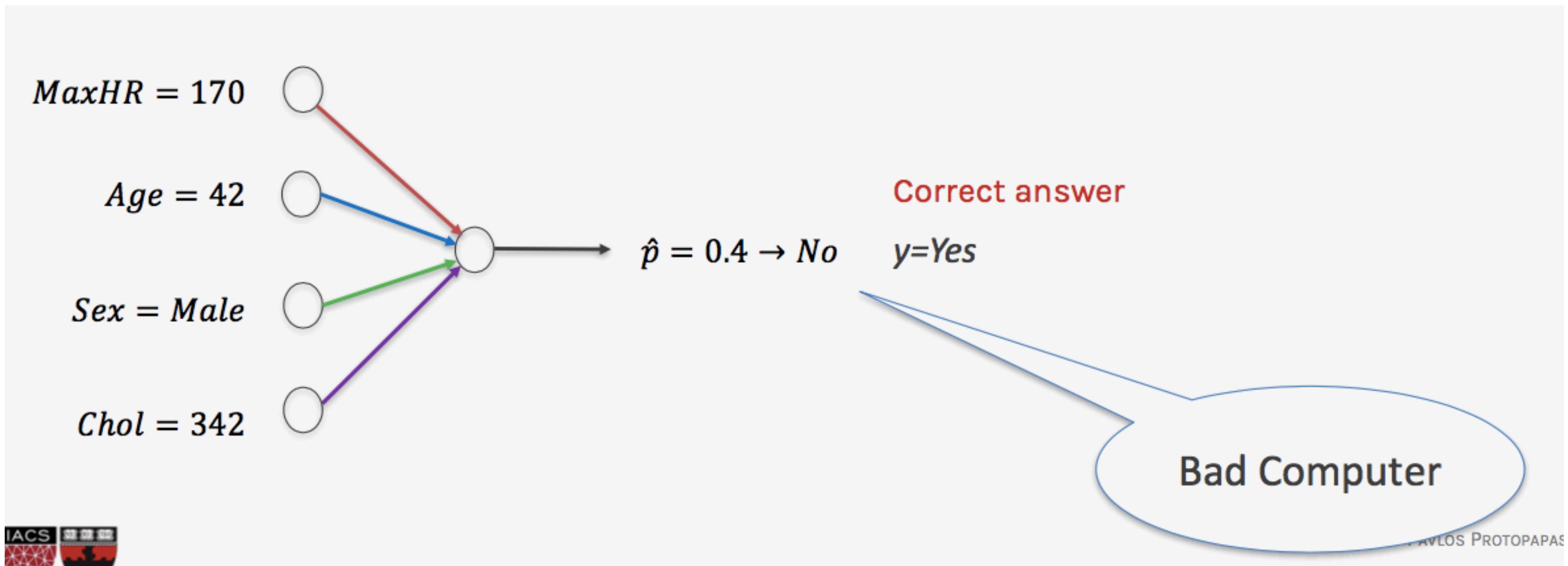
model.summary()
# Fit the model
model.fit(X_train, y_train, validation_data=(X_test, y_test),
           epochs=config.epochs,
           callbacks=[WandbCallback(data_type="image", labels=labels)])
```



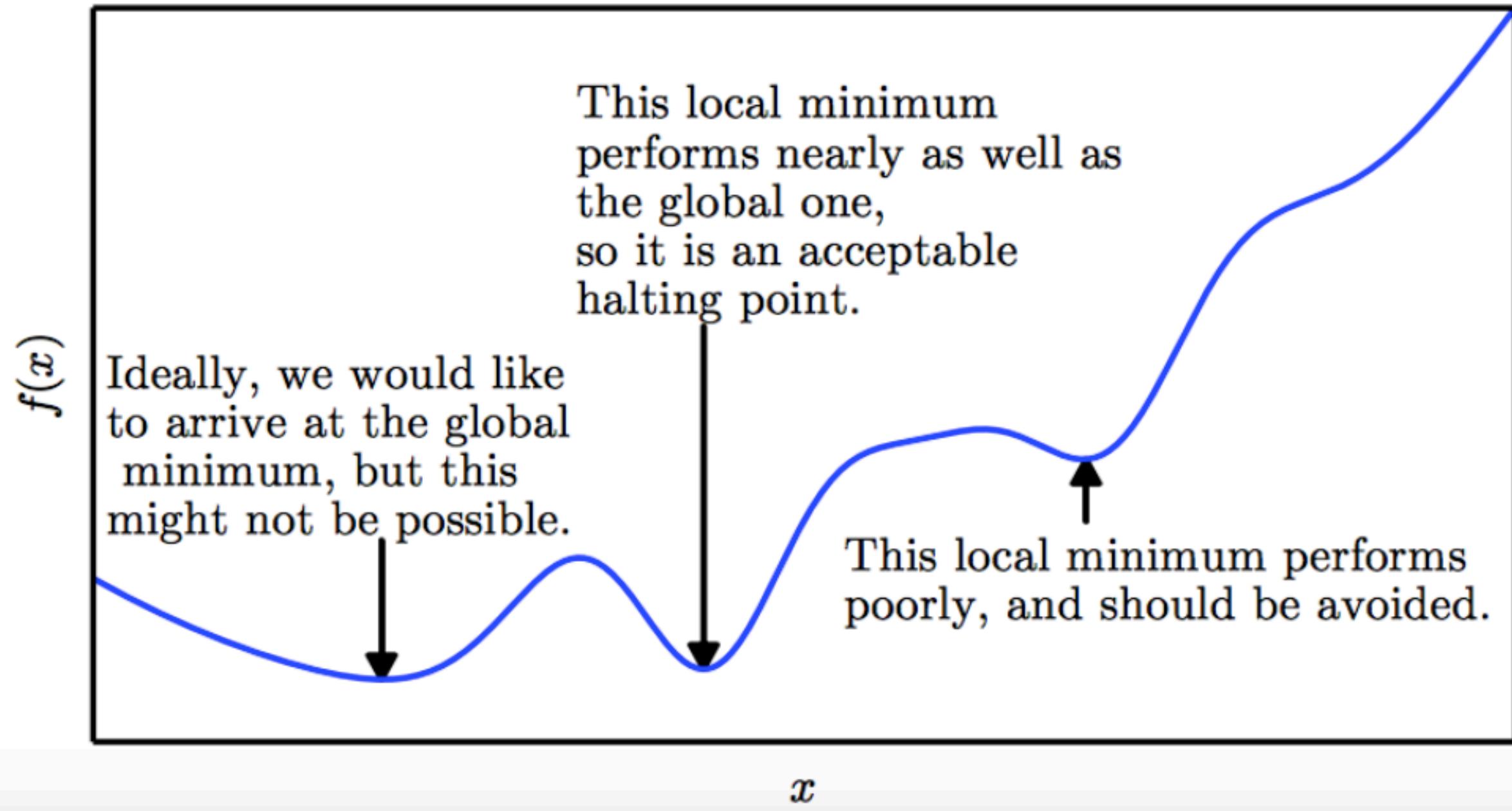




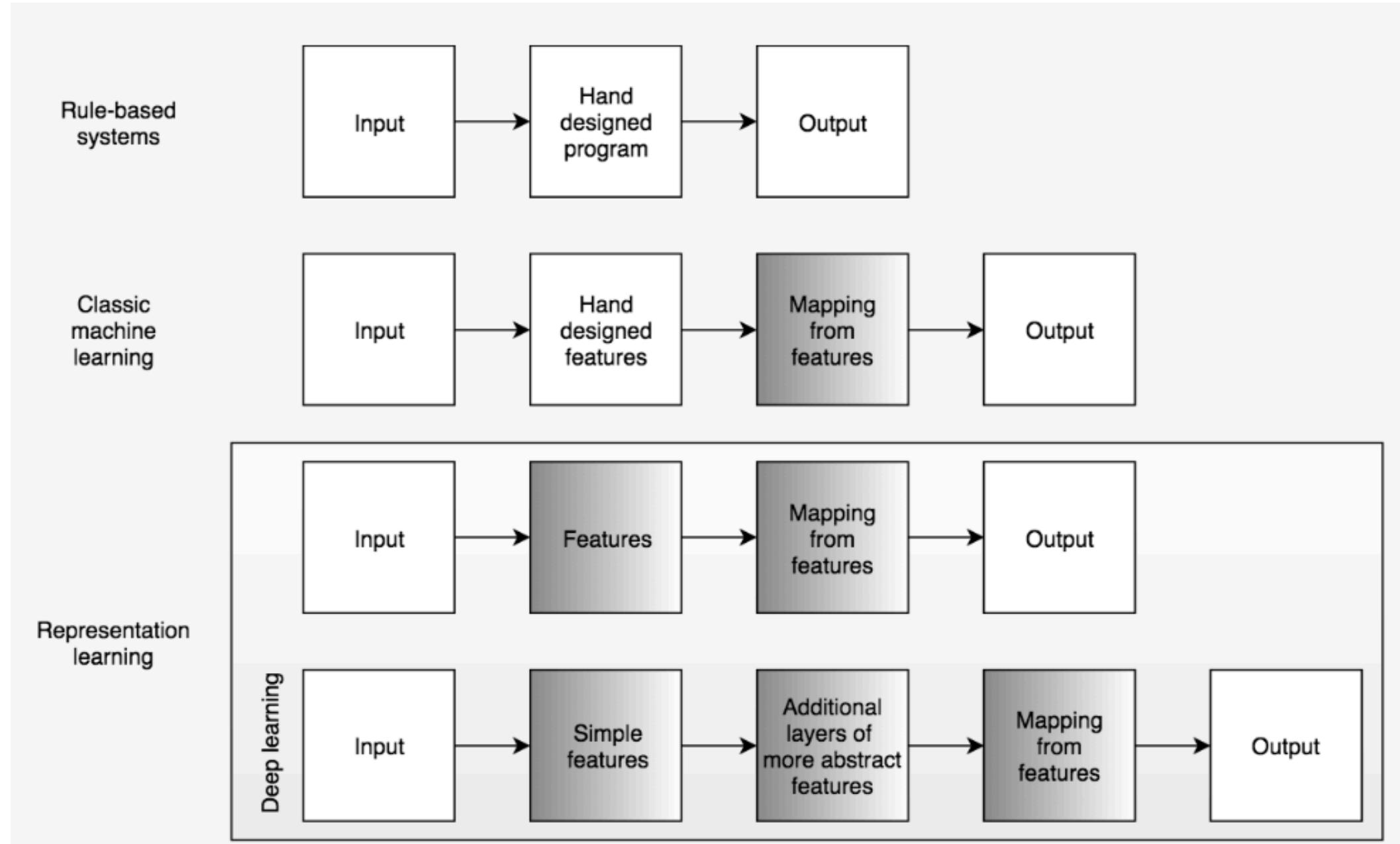
# We start with Garbage



# And get better, but we dont need to be best



# BIG IDEA: learn hierarchical representations



## And we need to help it

- MLPs are crude
- we must help computer
- thus use CNNs and RNNs which direct the representations
- directed representation learning is the basis for transfer learning

Three pillars

GPU

Automatic Differentiation

Representation Learning

# Statistics: Likelihoods

Output Type	Output Distribution	Output layer	Cost Function
Binary	Bernoulli	Sigmoid	Binary Cross Entropy
Discrete	Multinoulli	Softmax	Cross Entropy
Continuous	Gaussian	Linear	MSE
Continuous	Arbitrary	-	GANS

# Practicalities

## Problem: Gradient vanishing and explosion

Suppose  $\mathbf{W} = \begin{bmatrix} a & 0 \\ 0 & b \end{bmatrix}$ :

$$\begin{bmatrix} h_1^1 \\ h_2^1 \end{bmatrix} = \begin{bmatrix} a & 0 \\ 0 & b \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \quad \dots$$

$$\begin{bmatrix} h_1^n \\ h_2^n \end{bmatrix} = \begin{bmatrix} a^n & 0 \\ 0 & b^n \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

# Mitigations and Optimization

- Explosion: Gradient clipping
- Implosion: Skip Connections, Resnets, LSTM, Attention
- Curvature: Momentum, Gradient Accumulation
- Heterogeneity: adaptive learning rates

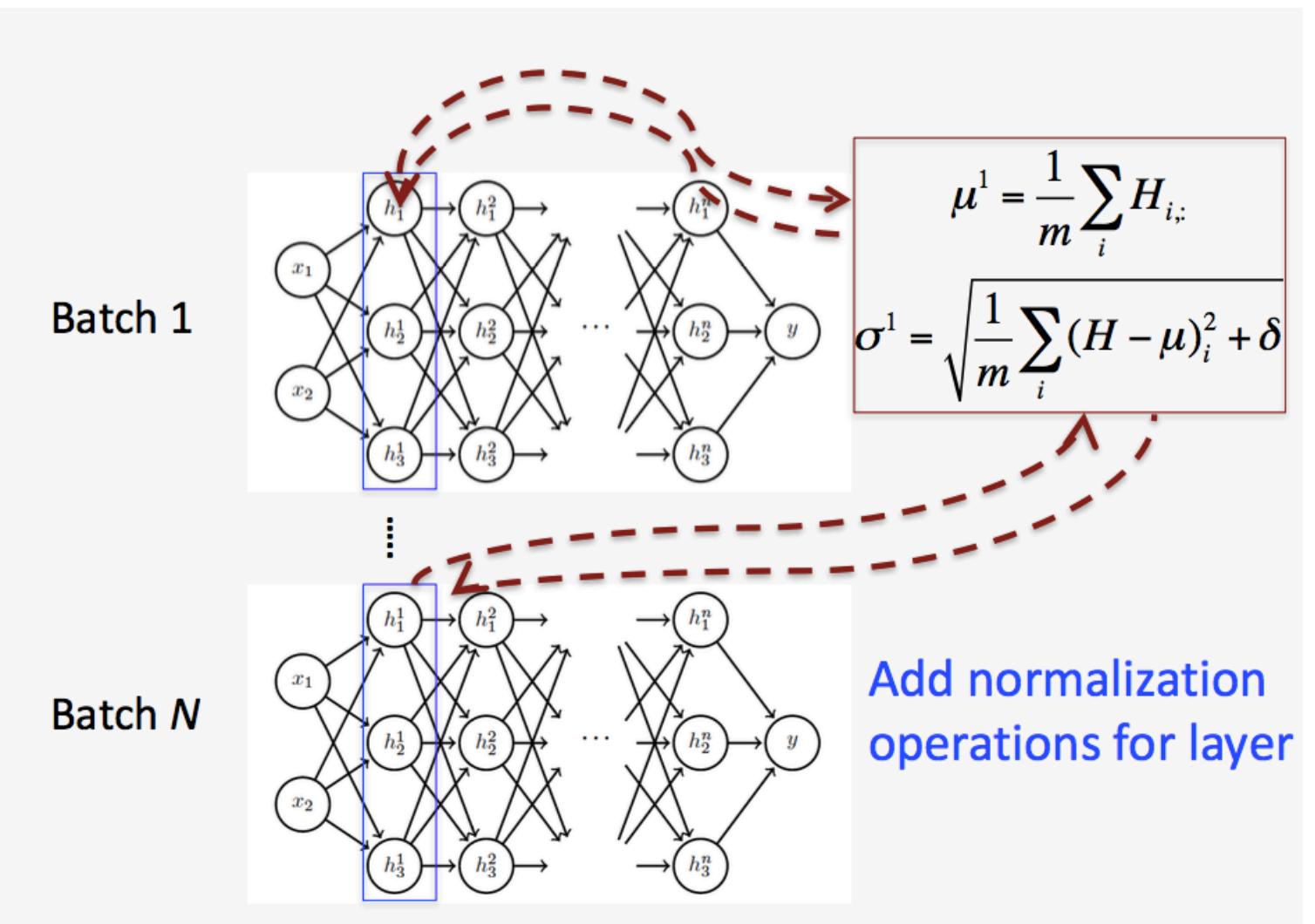
# Initialization

- Uniform or Normal: get unit variance
- Non saturated Bias for Relus
- Feature Normalization
- Batch Norm

# Batch Norm

$\gamma H' + \beta$

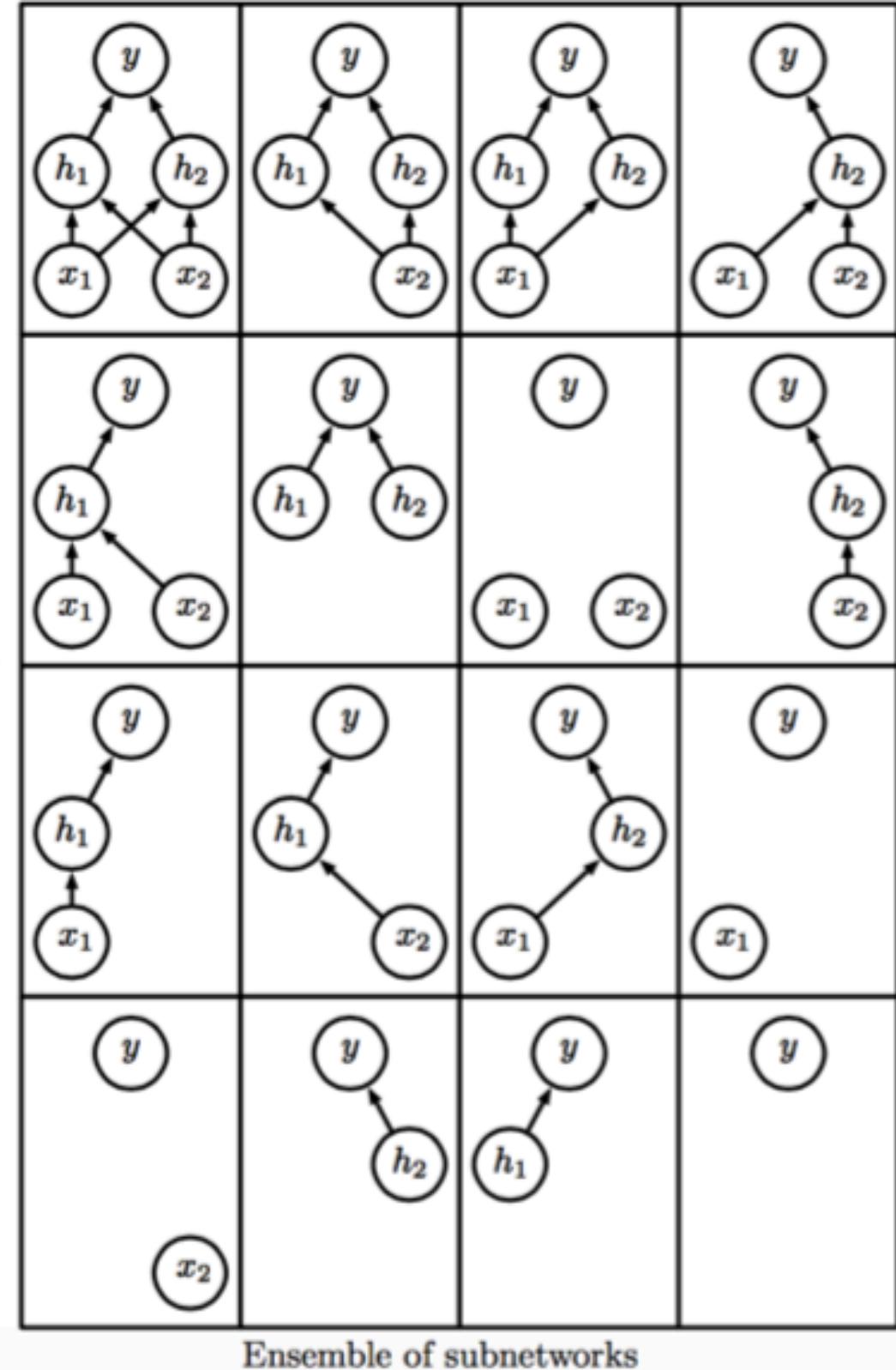
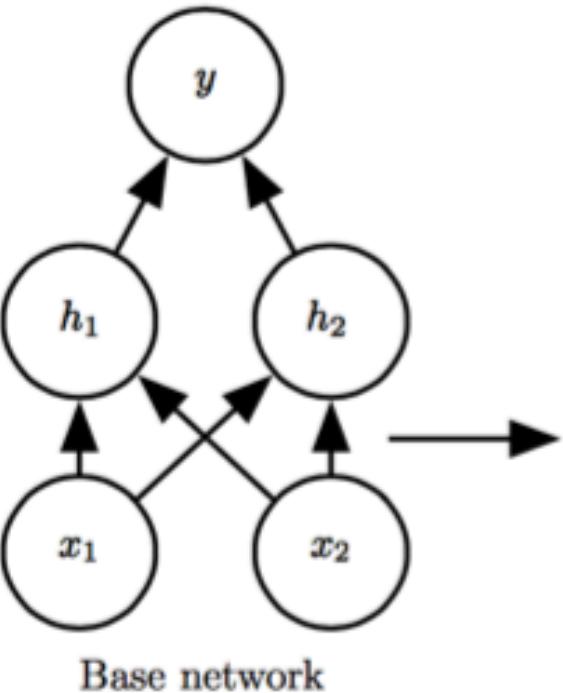
Learnable parameters



# First Overfit

# Dropout

```
model=Sequential()  
model.add(Flatten(input_shape=(img_width,img_height)))  
model.add(Dropout(config.dropout))  
model.add(Dense(config.hidden_nodes, activation='relu'))  
model.add(Dropout(config.dropout))  
model.add(Dense(num_classes, activation='softmax'))  
model.compile(loss='categorical_crossentropy', optimizer=config.optimizer,  
              metrics=['accuracy'])
```



# L1/L2:Weight Decay

We used to optimize:

Change to ...

$$W^{(i+1)} = W^{(i)} - \lambda \frac{\partial J}{\partial W}$$
$$J_R(W; X, y) = J(W; X, y) + \frac{1}{2} \alpha W^2$$
$$W^{(i+1)} = W^{(i)} - \lambda \frac{\partial J}{\partial W} \boxed{- \lambda \alpha W}$$

weights decay in proportion to its size.

Biases not penalized

$L_2$  regularization:

- Decay of weights
- MAP estimation with Gaussian prior

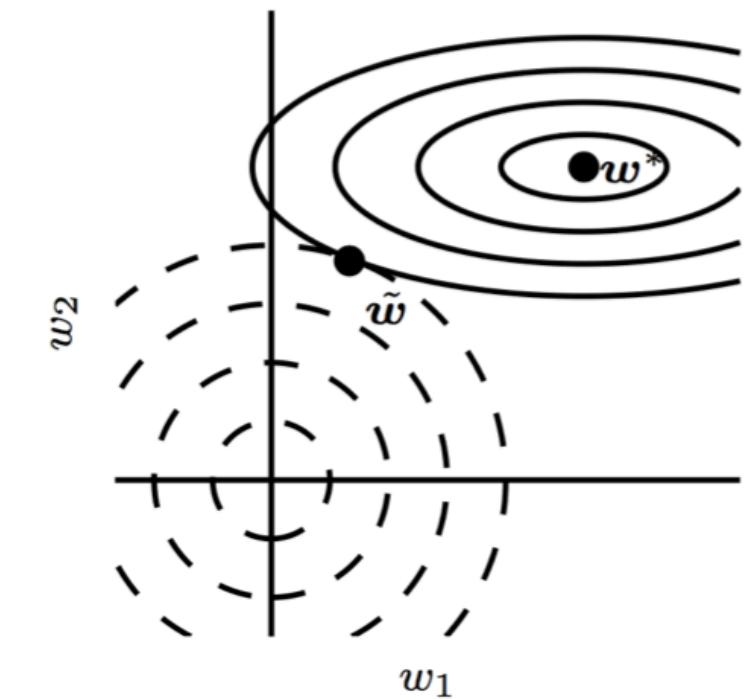
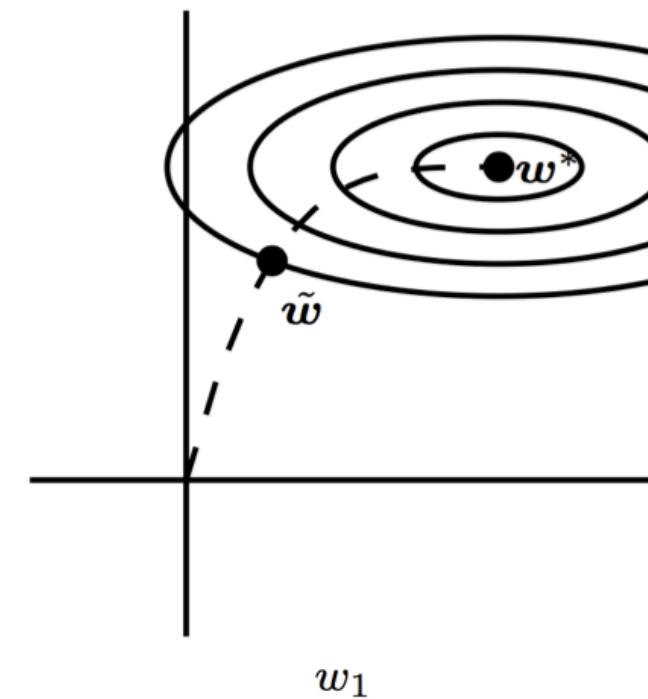
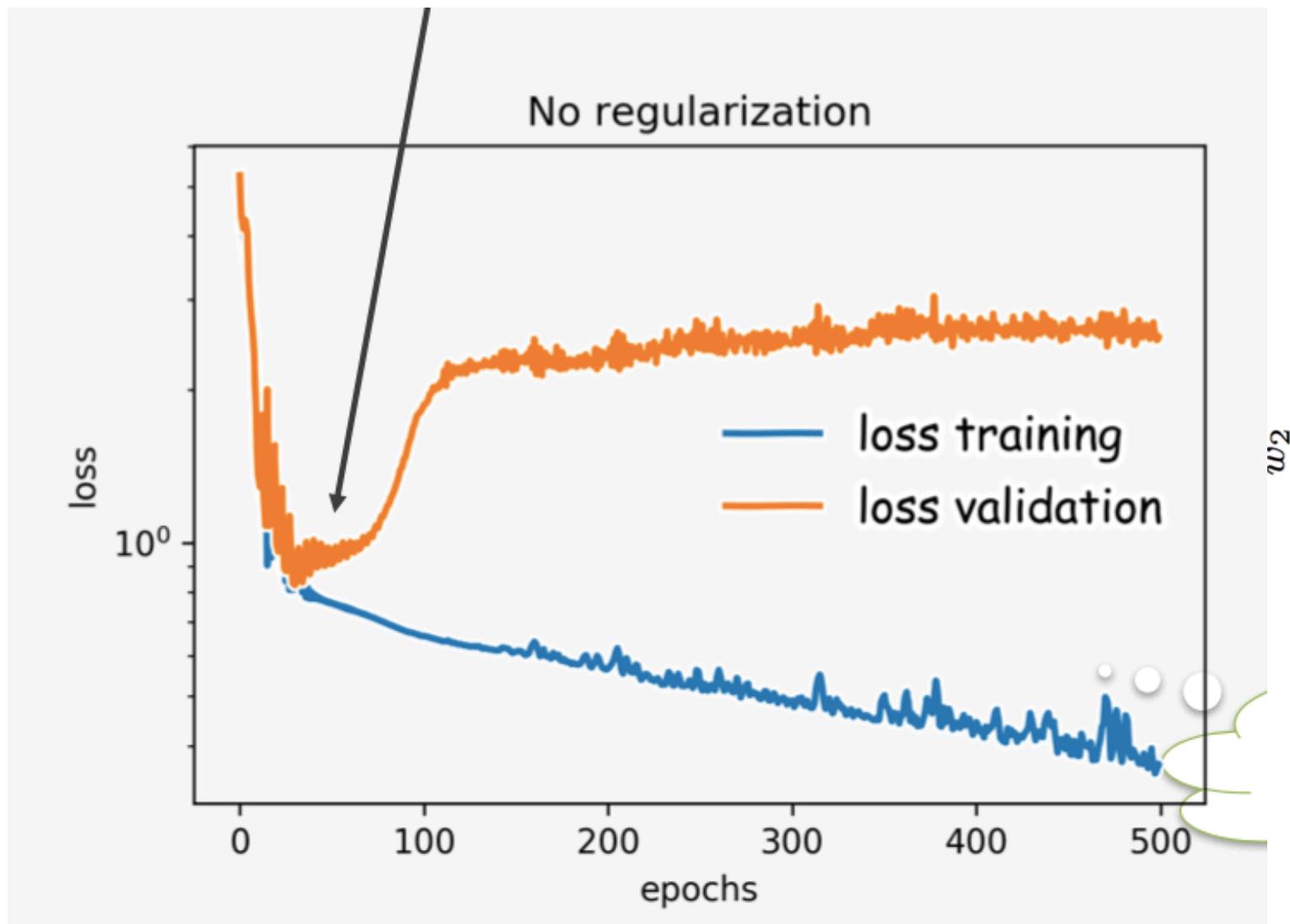
$$\Omega(W) = \frac{1}{2} \| W \|_2^2$$

$L_1$  regularization:

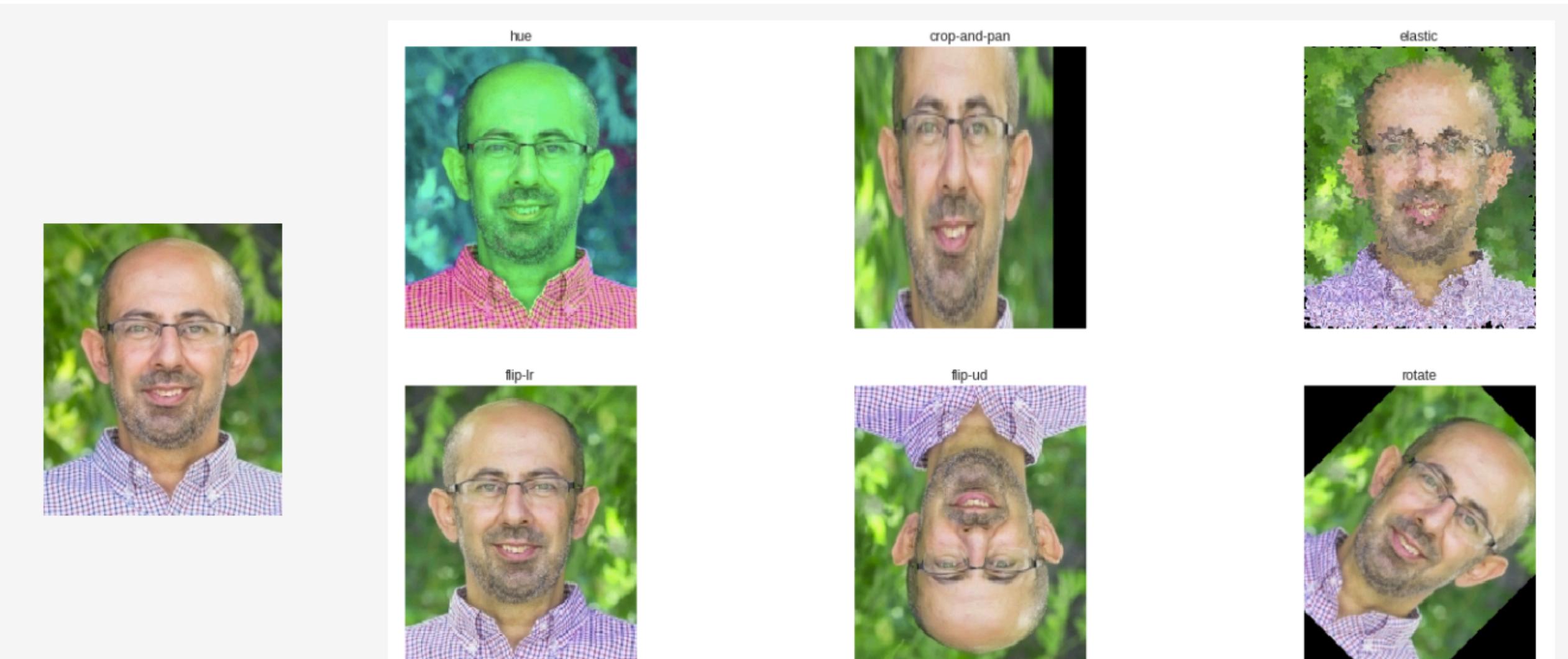
- encourages sparsity
- MAP estimation with Laplacian prior

$$\Omega(W) = \frac{1}{2} \| W \|_1$$

# Early Stopping



# Data Aug

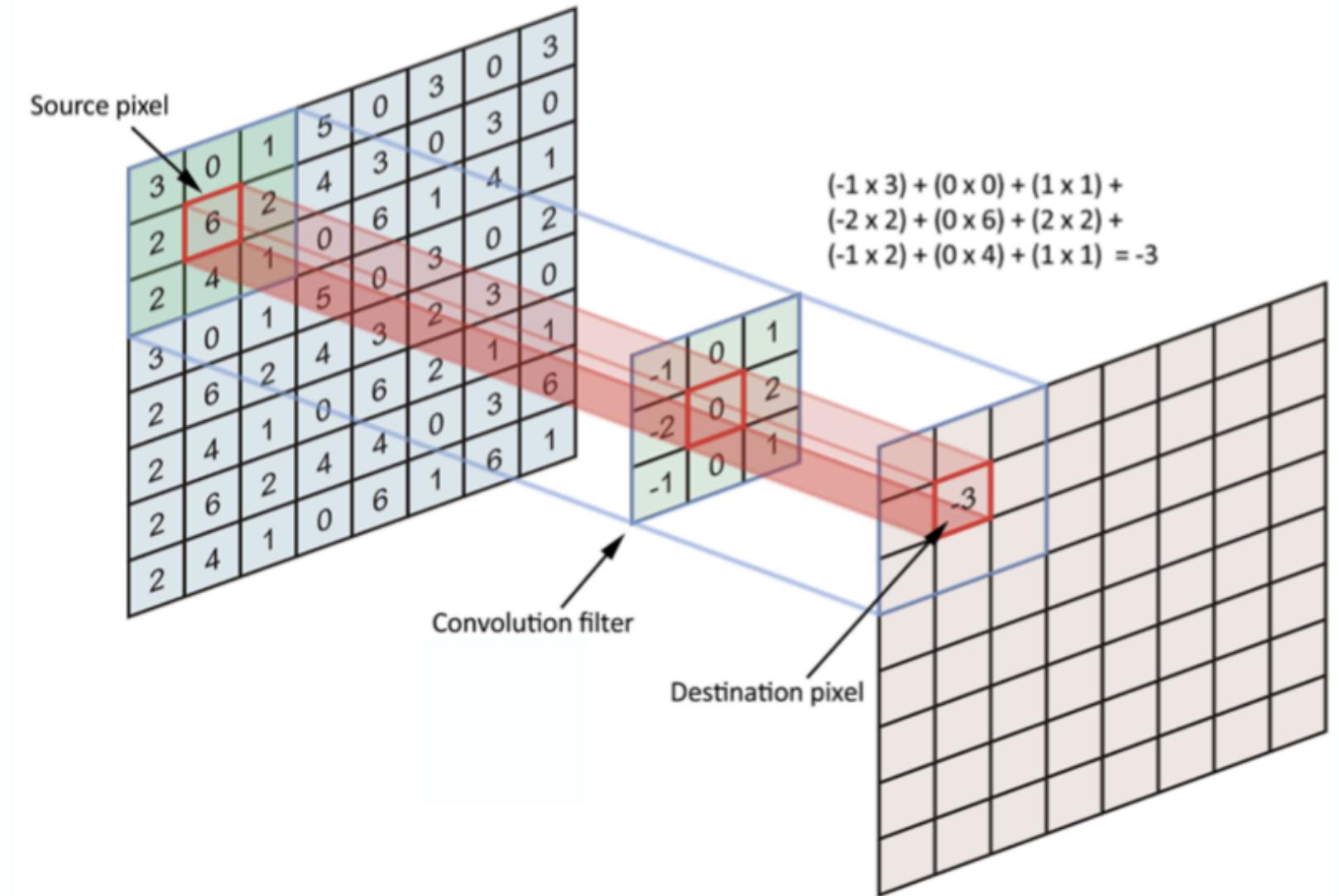
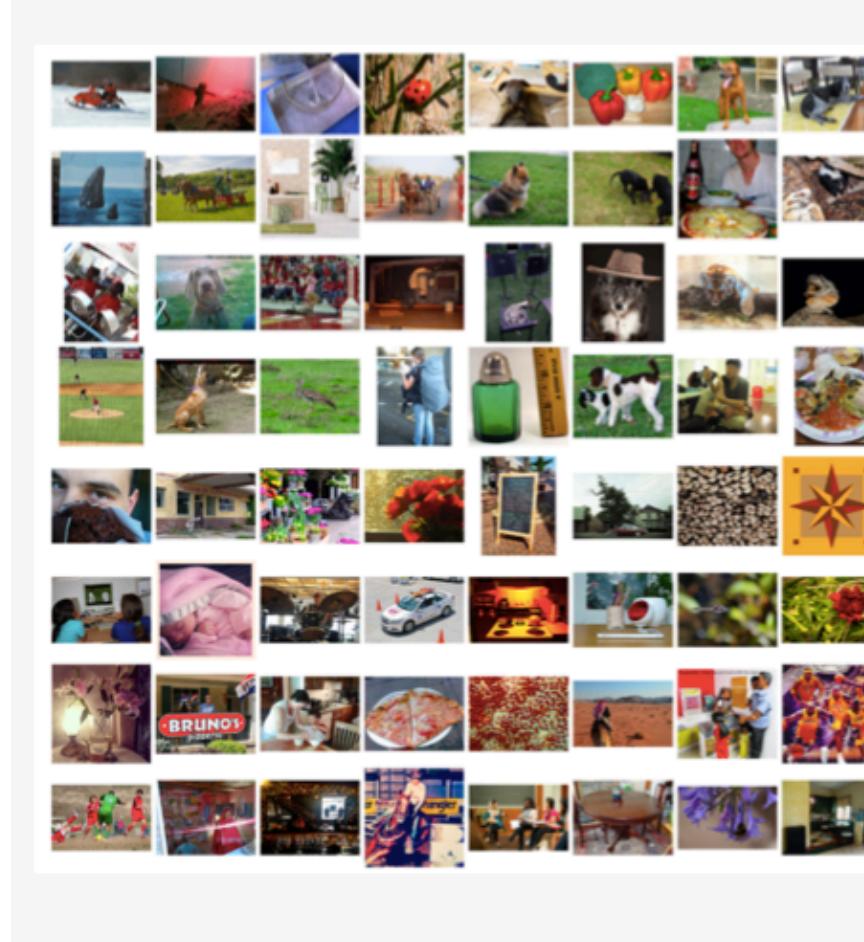


## Others

- Bagging
- Adversarial Examples
- Sparse Activations (on activations not weights)

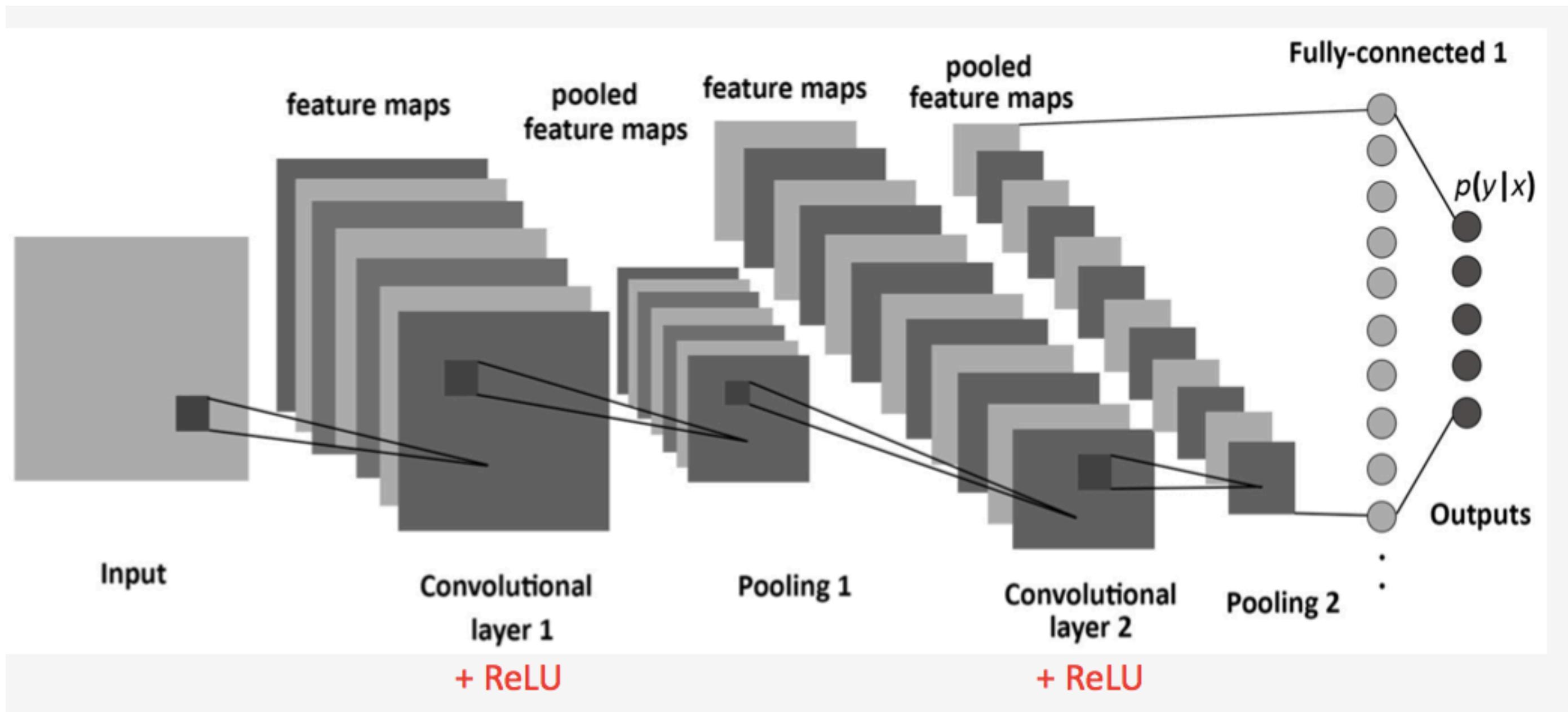
# Convolutional Neural Networks

# Exploding params: do weight tying

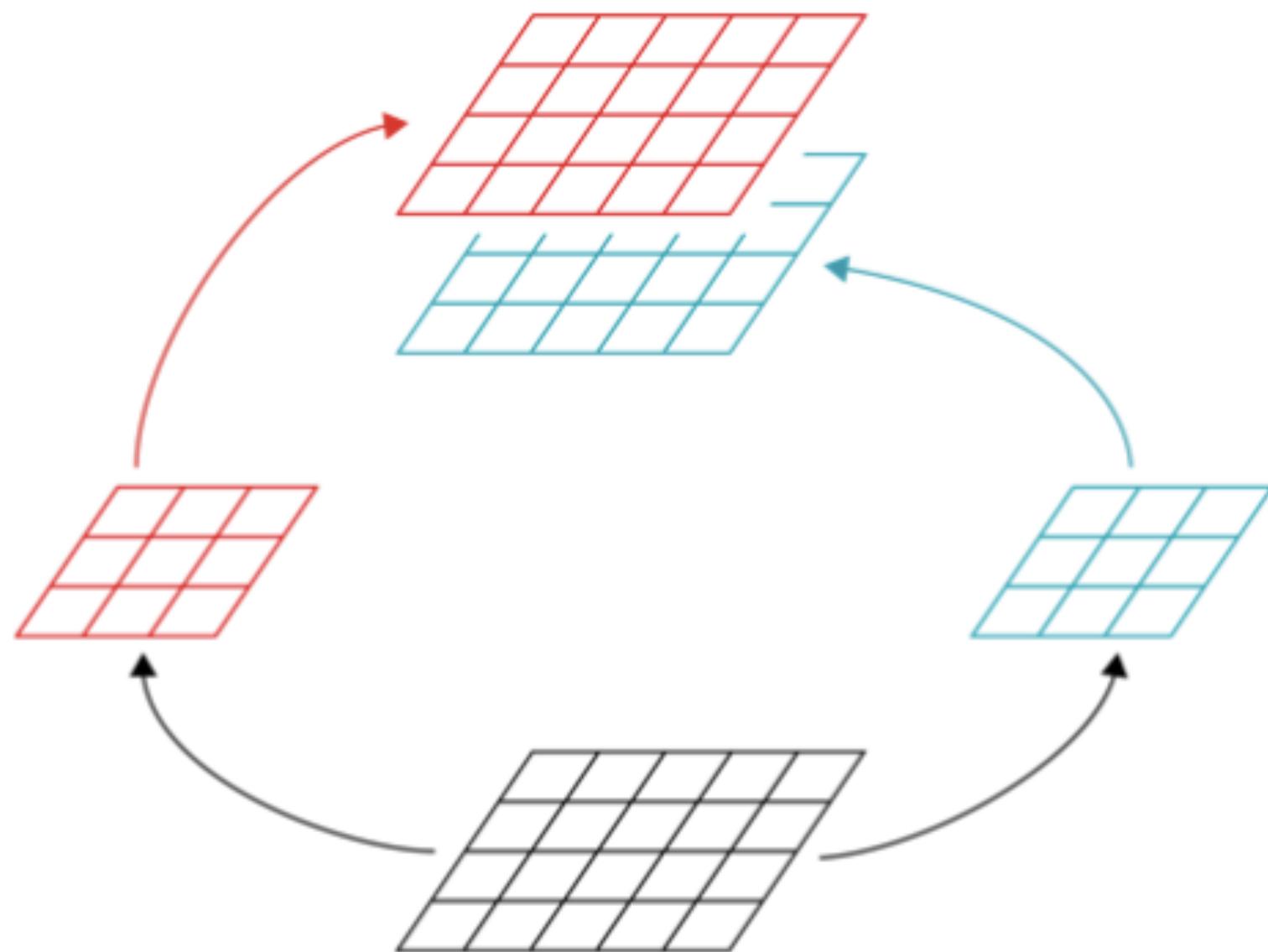


- 3073 params cifar-10
- 150129 Imagenet
- Nearer pixels are related
- And there are symmetries

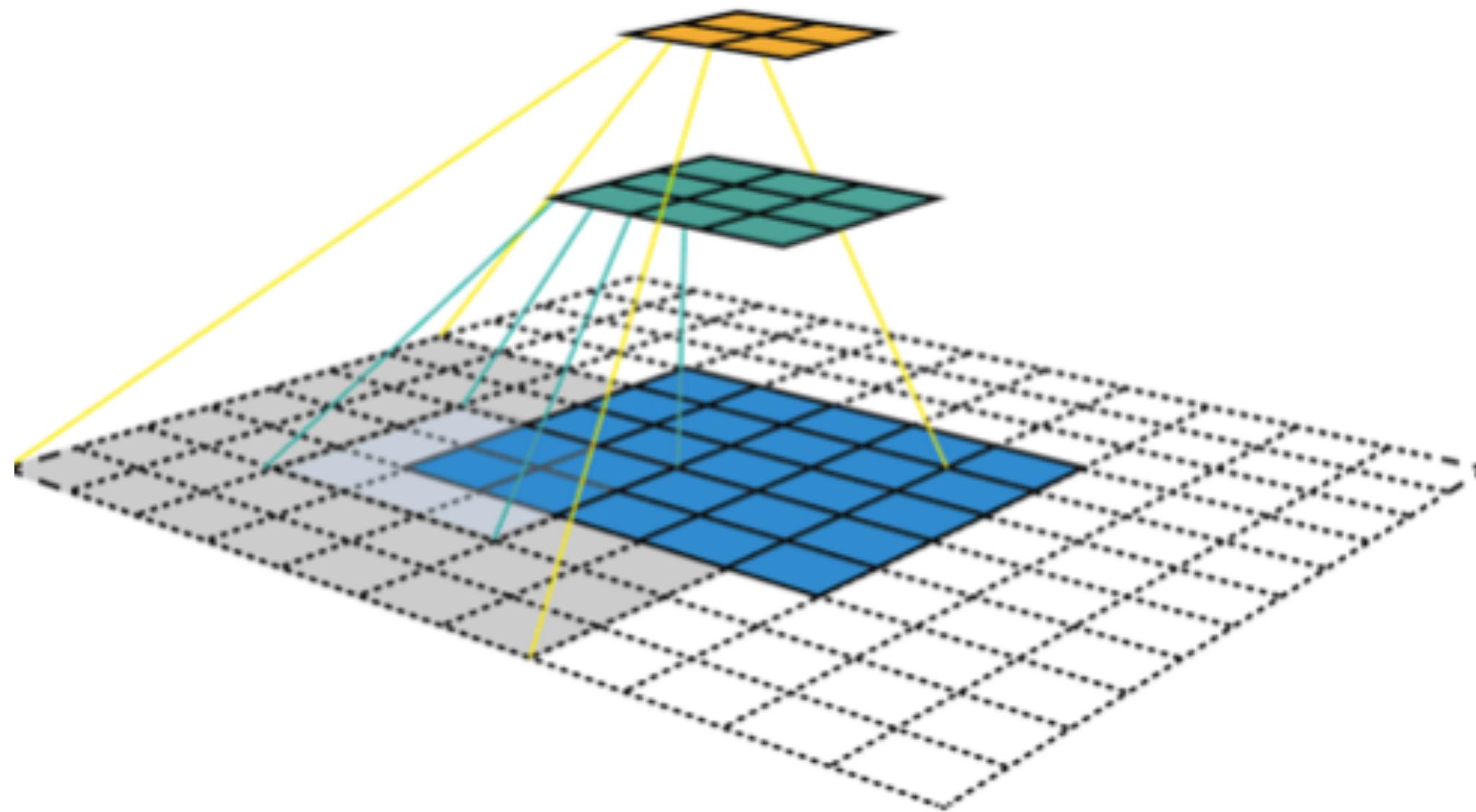
# Architecture



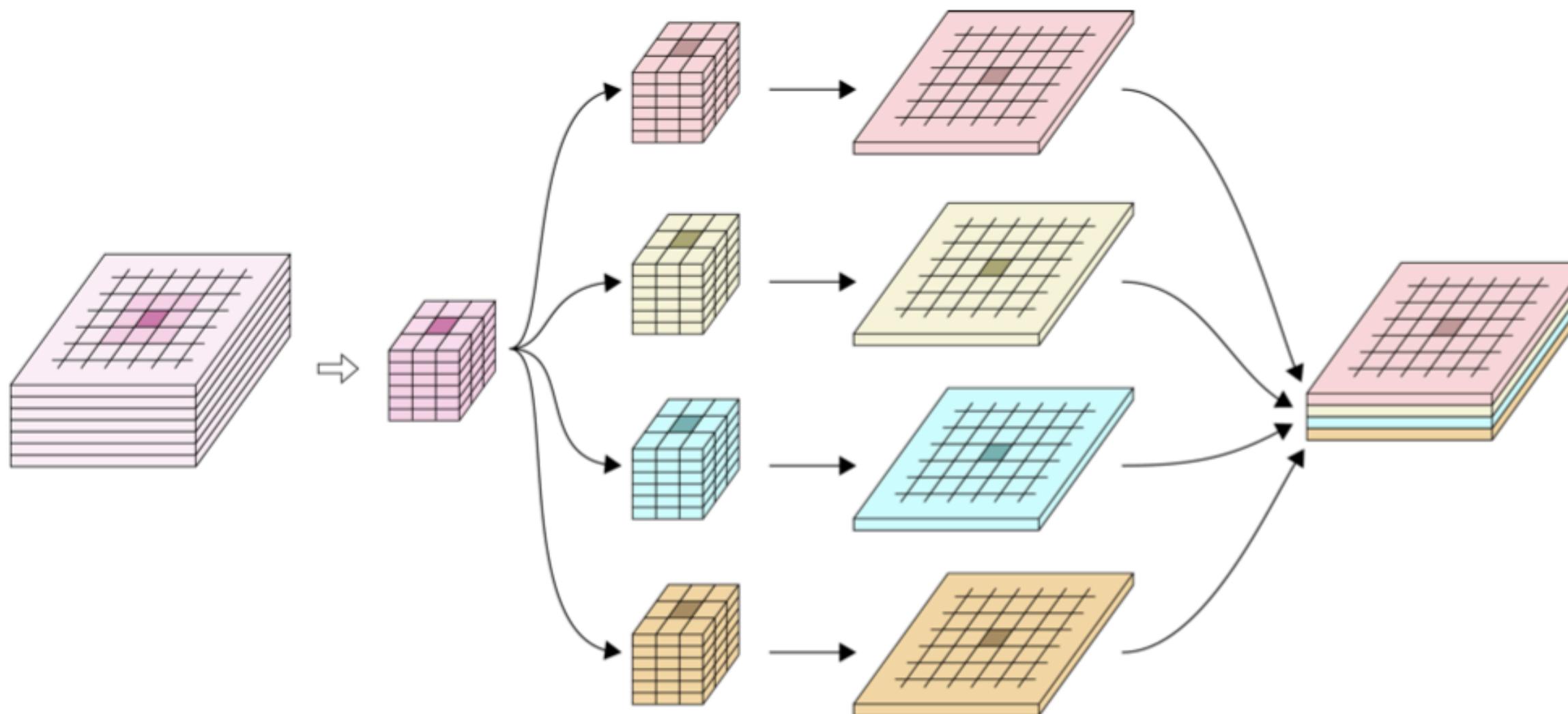
# Featuremaps



# Receptive Field



# How featuremaps in filters work



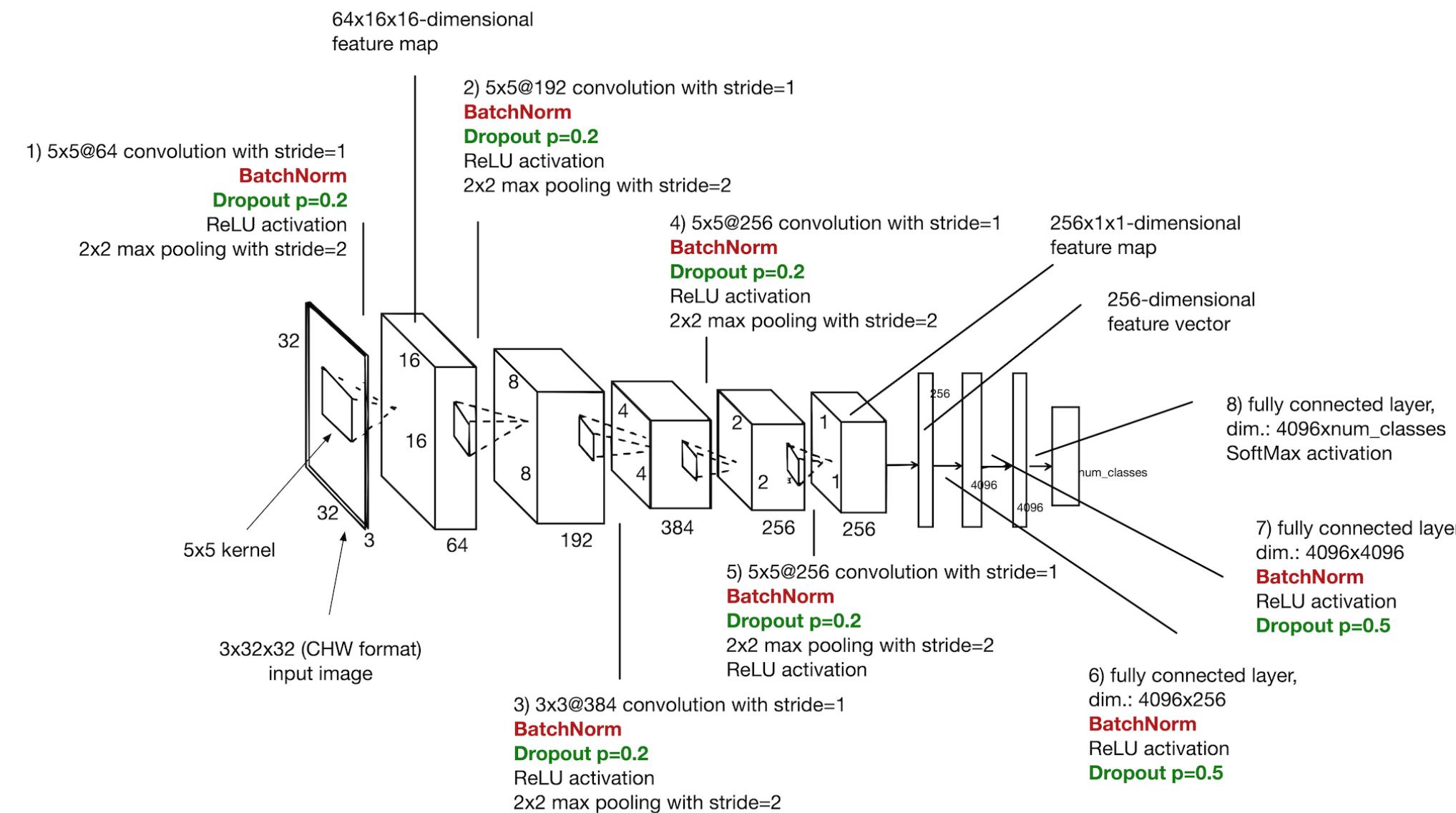
# Things to do

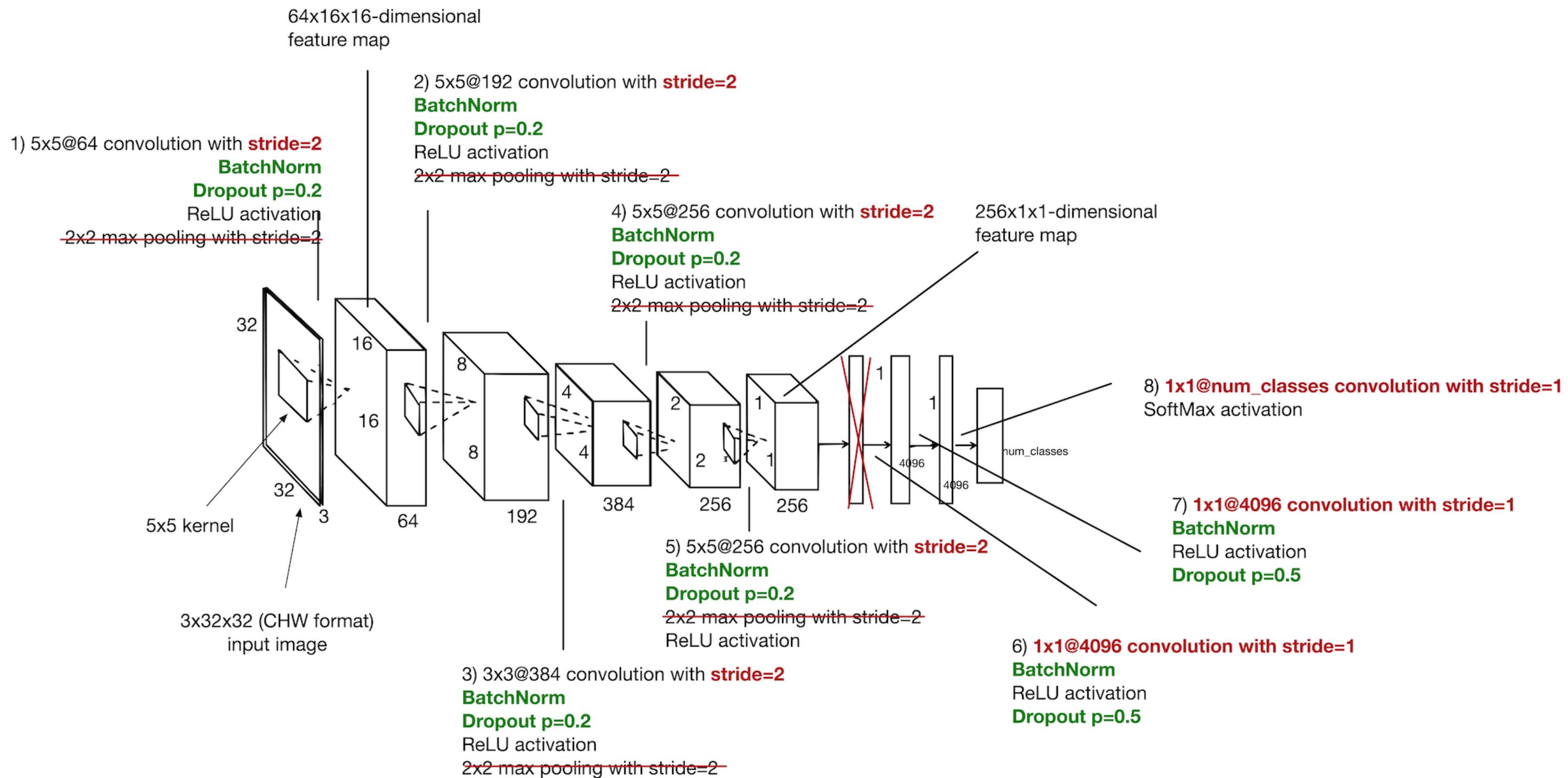
- striding
- padding
- pooling
- upsampling
- 1x1 convolutions
- globalavgpooling
- reshaping, Densing

# Recursive Learning

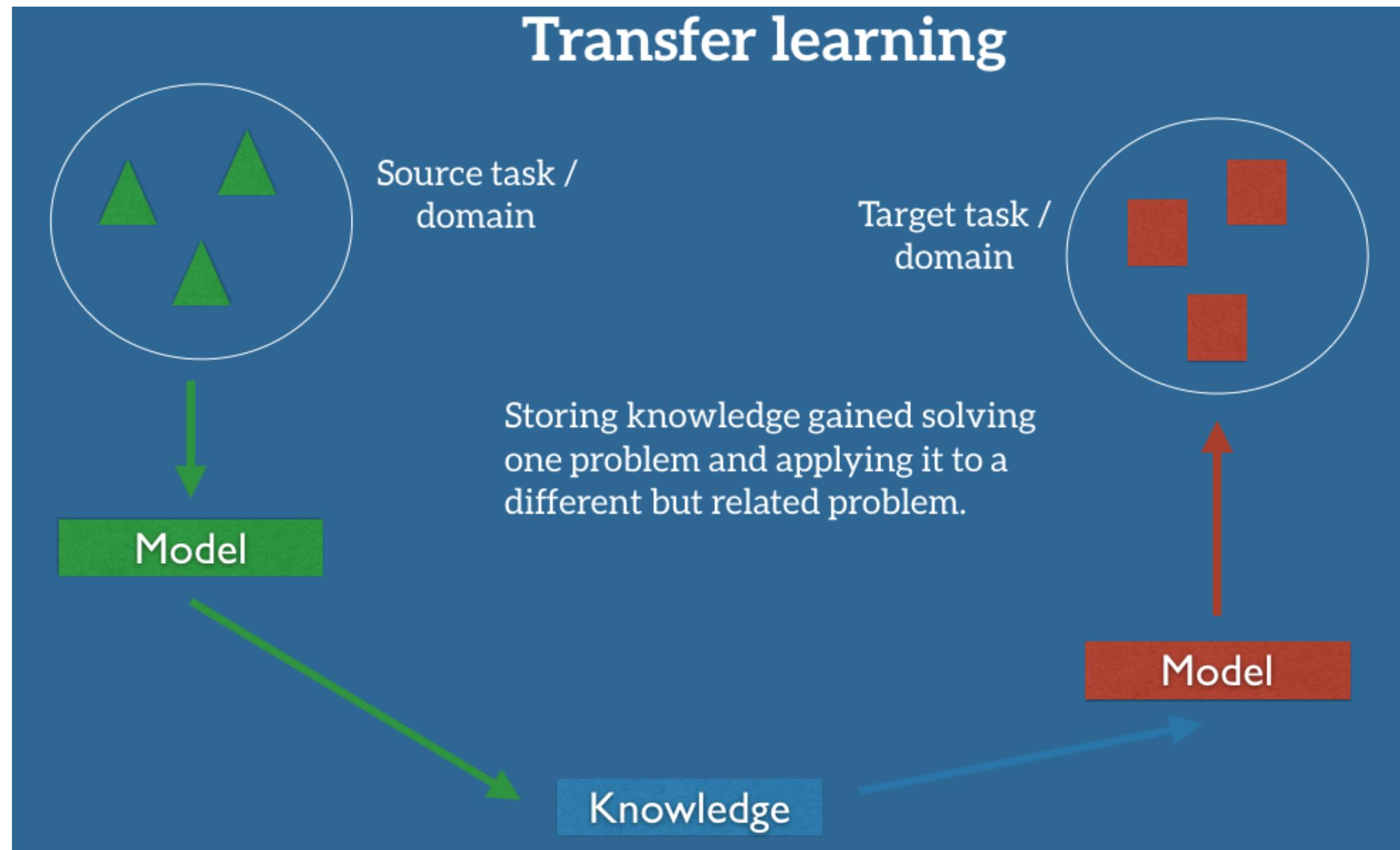


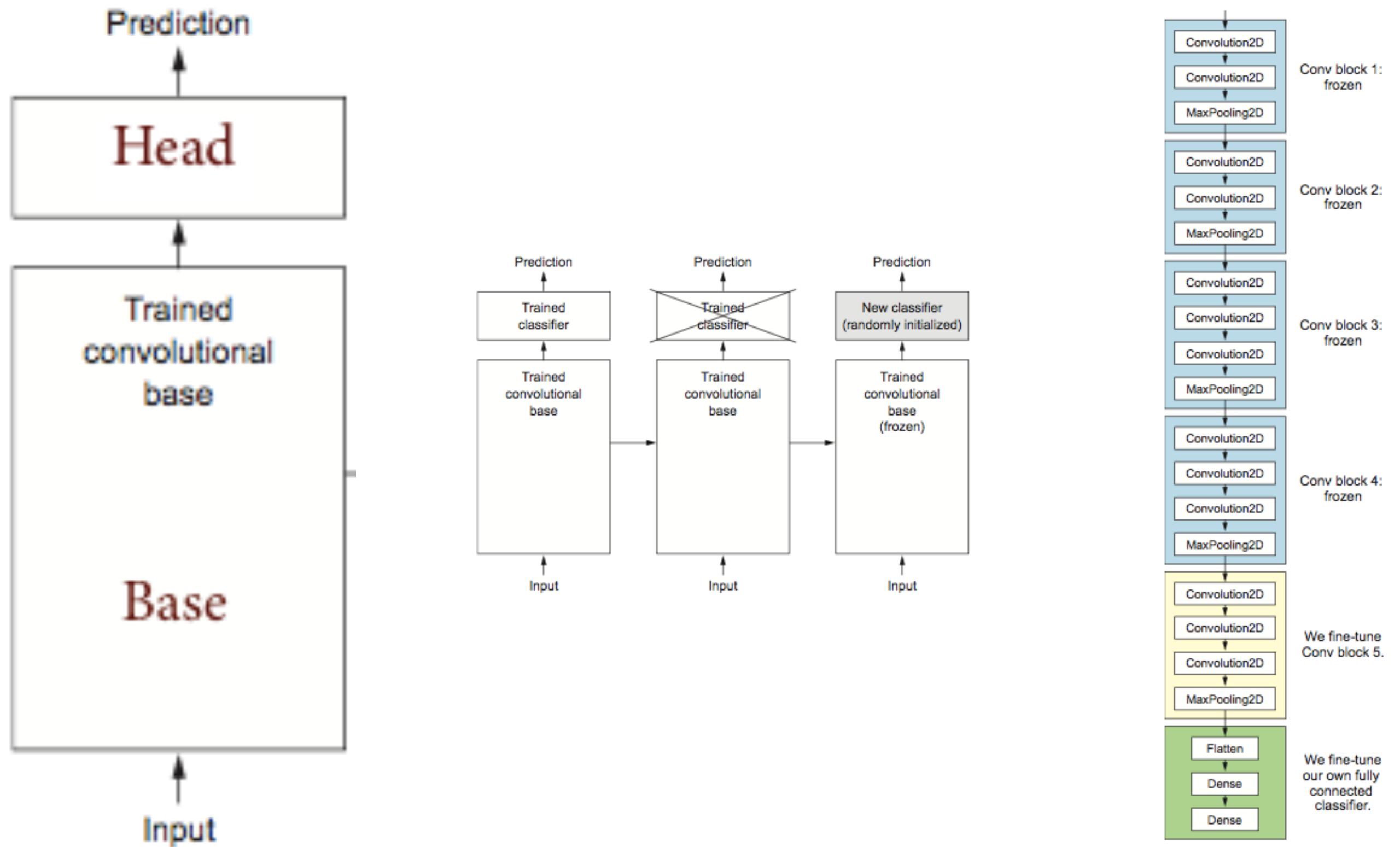
# Architectures





# Transfer Learning





Language  
Modeling  
+ Embeddings

# Basic

```
tokenizer = text.Tokenizer(num_words=config.vocab_size)
tokenizer.fit_on_texts(X_train)
X_train = tokenizer.texts_to_matrix(X_train)
X_test = tokenizer.texts_to_matrix(X_test)

bow_model = LogisticRegression()
bow_model.fit(X_train, y_train)

pred_train = bow_model.predict(X_train)
acc = np.sum(pred_train==y_train)/len(pred_train)

pred_test = bow_model.predict(X_test)
val_acc = np.sum(pred_test==y_test)/len(pred_test)
```

# Language Modeling

Thou shalt not make **a machine in** the likeness of a human mind

Sliding window across running text

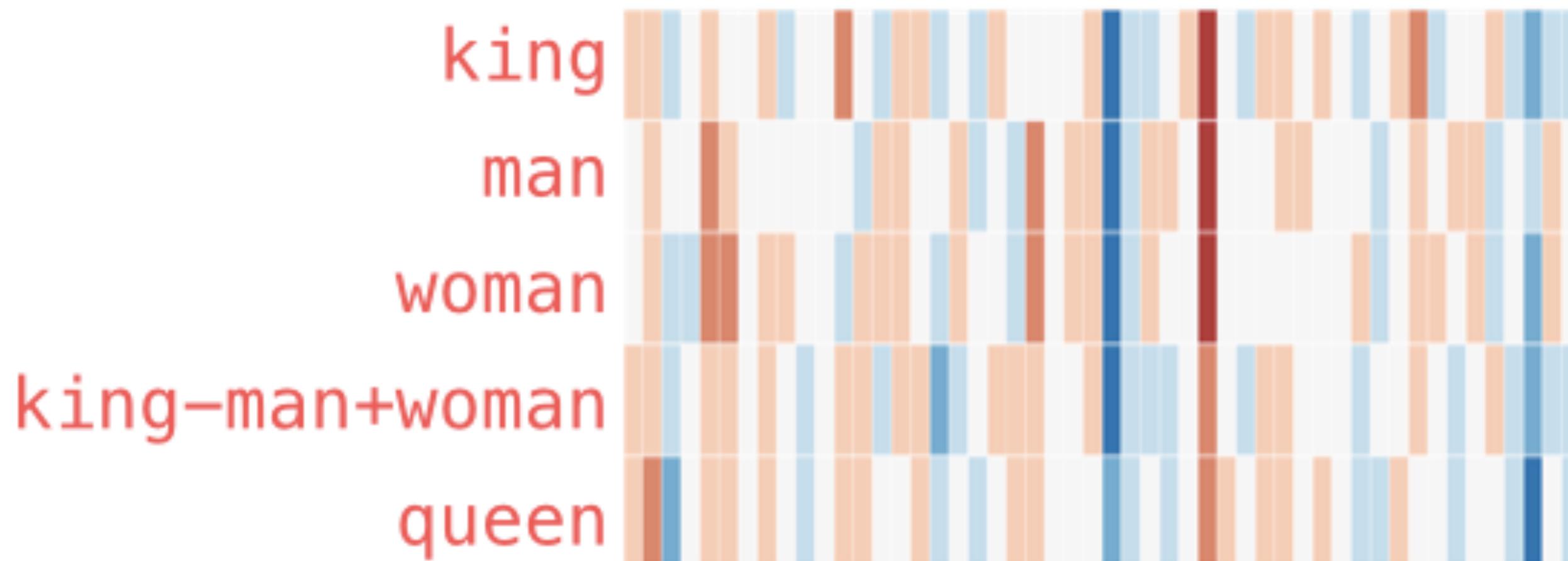
thou	shalt	not	make	a	machine	in	the	...
thou	shalt	not	make	a	machine	in	the	
thou	shalt	not	make	a	machine	in	the	
thou	shalt	not	make	a	machine	in	the	
thou	shalt	not	make	a	machine	in	the	

Dataset

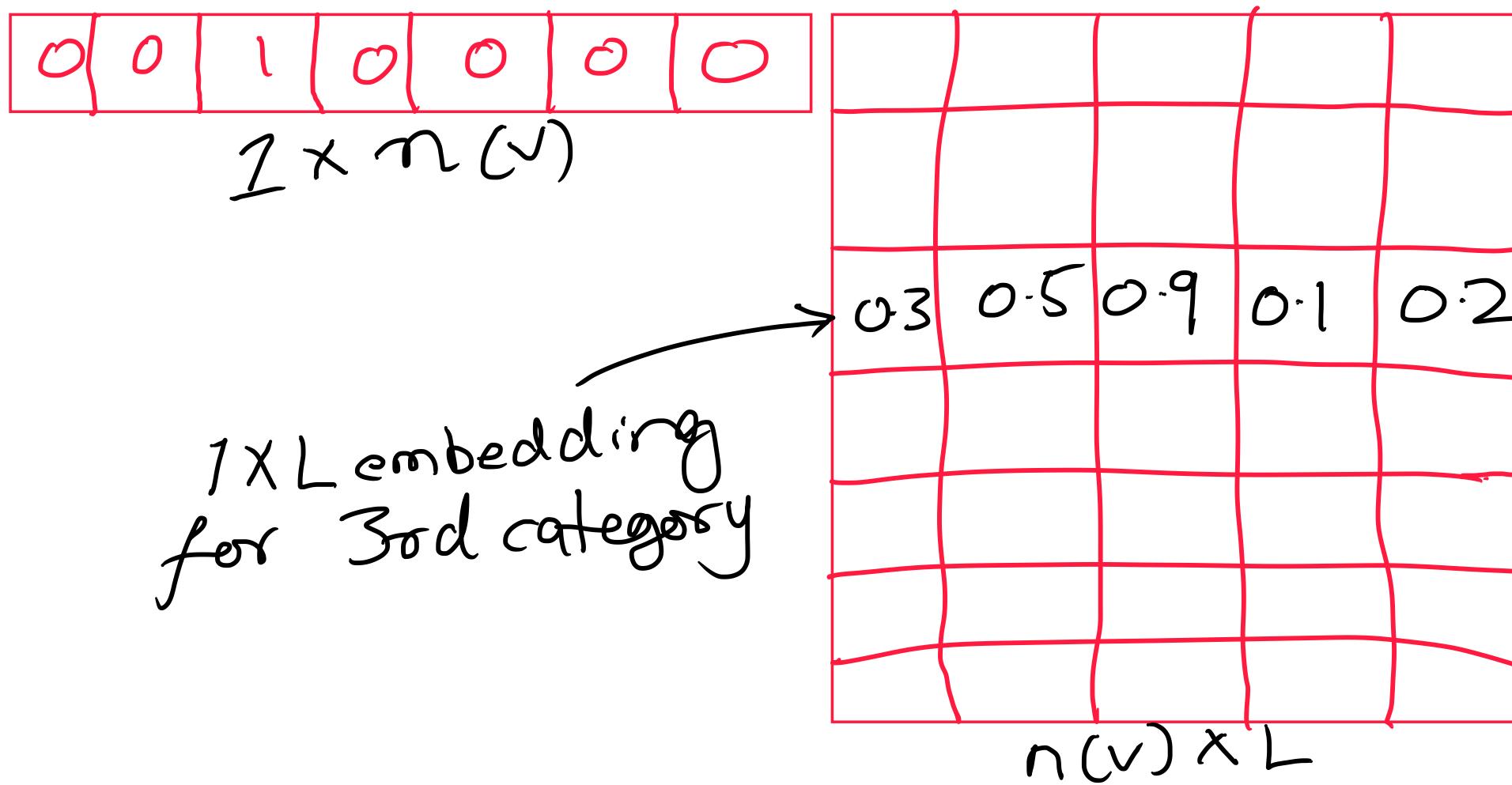
input 1	input 2	output
thou	shalt	not
shalt	not	make
not	make	a
make	a	machine
a	machine	in

# Word Embeddings

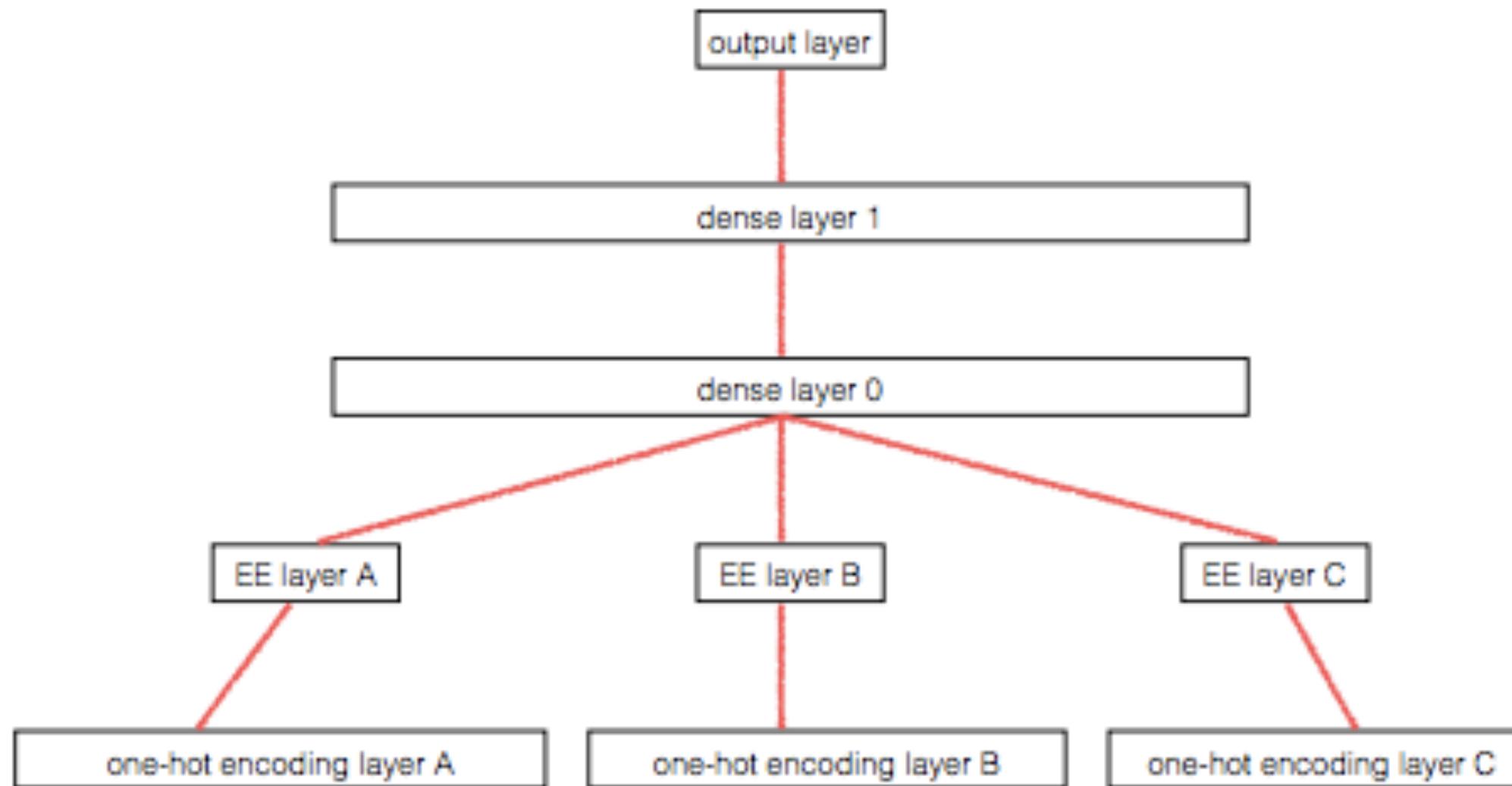
$$\text{king} - \text{man} + \text{woman} \approx \text{queen}$$



# Embeddings are Linear Regression



# Learn Embeddings along with Task at hand



# Application: Recommendations

$$Y_{35}^{baseline} = \mu + \bar{\theta} \cdot I_3 + \bar{\gamma} \cdot I_5 = \mu + [I_3, I_5] \cdot [\bar{\theta}, \bar{\gamma}]$$

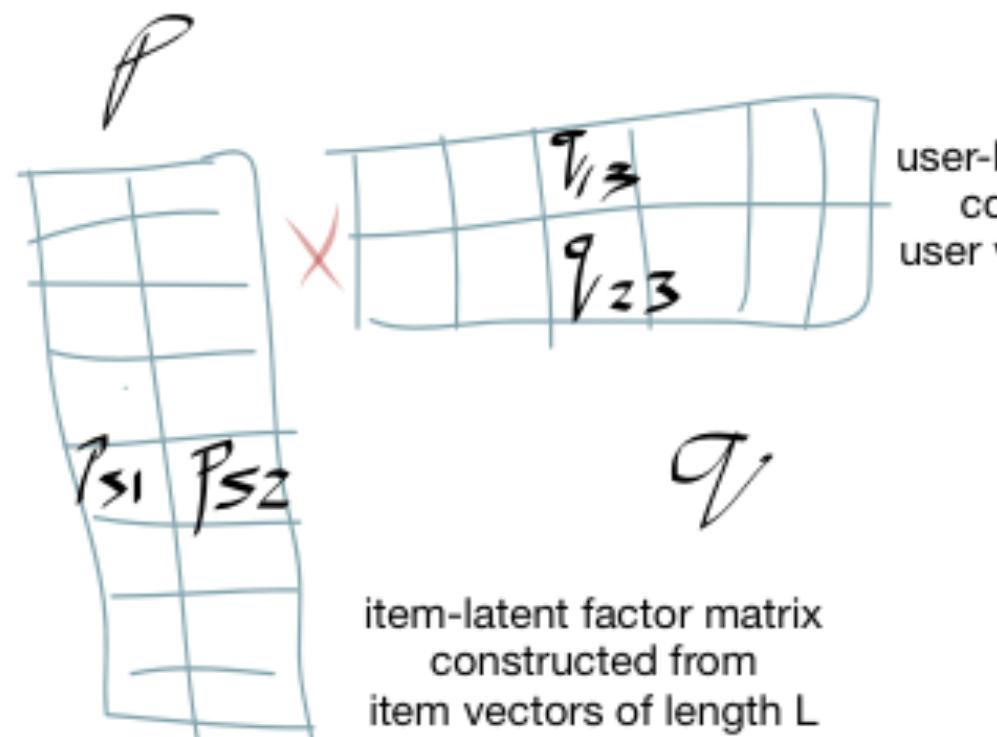
$$Y_{35} = \mu \begin{pmatrix} 1 \\ 1 \\ 1 \\ 1 \\ 1 \end{pmatrix} + \begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} \theta_3 \\ \theta_4 \\ \theta_5 \\ \gamma_1 \\ \gamma_2 \\ \gamma_3 \\ \gamma_4 \\ \gamma_5 \end{pmatrix}$$

Diagram illustrating the matrix equation for recommendations:

- Y<sub>35</sub>**: N x 1 vector of ratings.
- μ**: N x 1 vector of intercepts.
- Design (Feature) Matrix**: N x (U+M) matrix where columns represent users (U) and items (M). The matrix has 8 rows (N) and 9 columns (U+M). The first column is circled in green and labeled "Intercept". The next 8 columns are circled in red and labeled "(1 X U)" and "(1 X M)" respectively.
- Coefficients**: (U+M) x 1 vector containing user coefficients  $\theta_3, \theta_4, \theta_5$  and item coefficients  $\gamma_1, \gamma_2, \gamma_3, \gamma_4, \gamma_5$ .

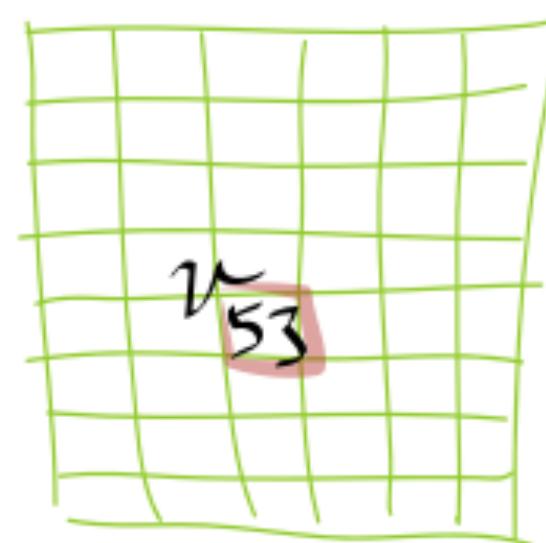
$$Y_{35}^{baseline} = \mu + \theta_3 + \gamma_5$$

L=2 latent factors



user-latent factor matrix  
constructed from  
user vectors of length L

item-latent factor matrix  
constructed from  
item vectors of length L



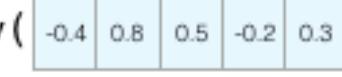
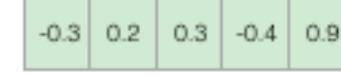
The product gives  
the user-item  
residual ranking matrix

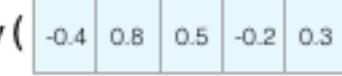
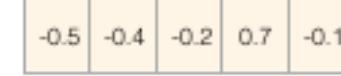
# The reasons for recommendations: similarity and FP

	Trait #1	Trait #2	Trait #3	Trait #4	Trait #5
Jay	-0.4	0.8	0.5	-0.2	0.3

Person #1	-0.3	0.2	0.3	-0.4	0.9
-----------	------	-----	-----	------	-----

Person #2	-0.5	-0.4	-0.2	0.7	-0.1
-----------	------	------	------	-----	------

`cosine_similarity(`  ,  ) = 0.66 ✓

`cosine_similarity(`  ,  ) = -0.37

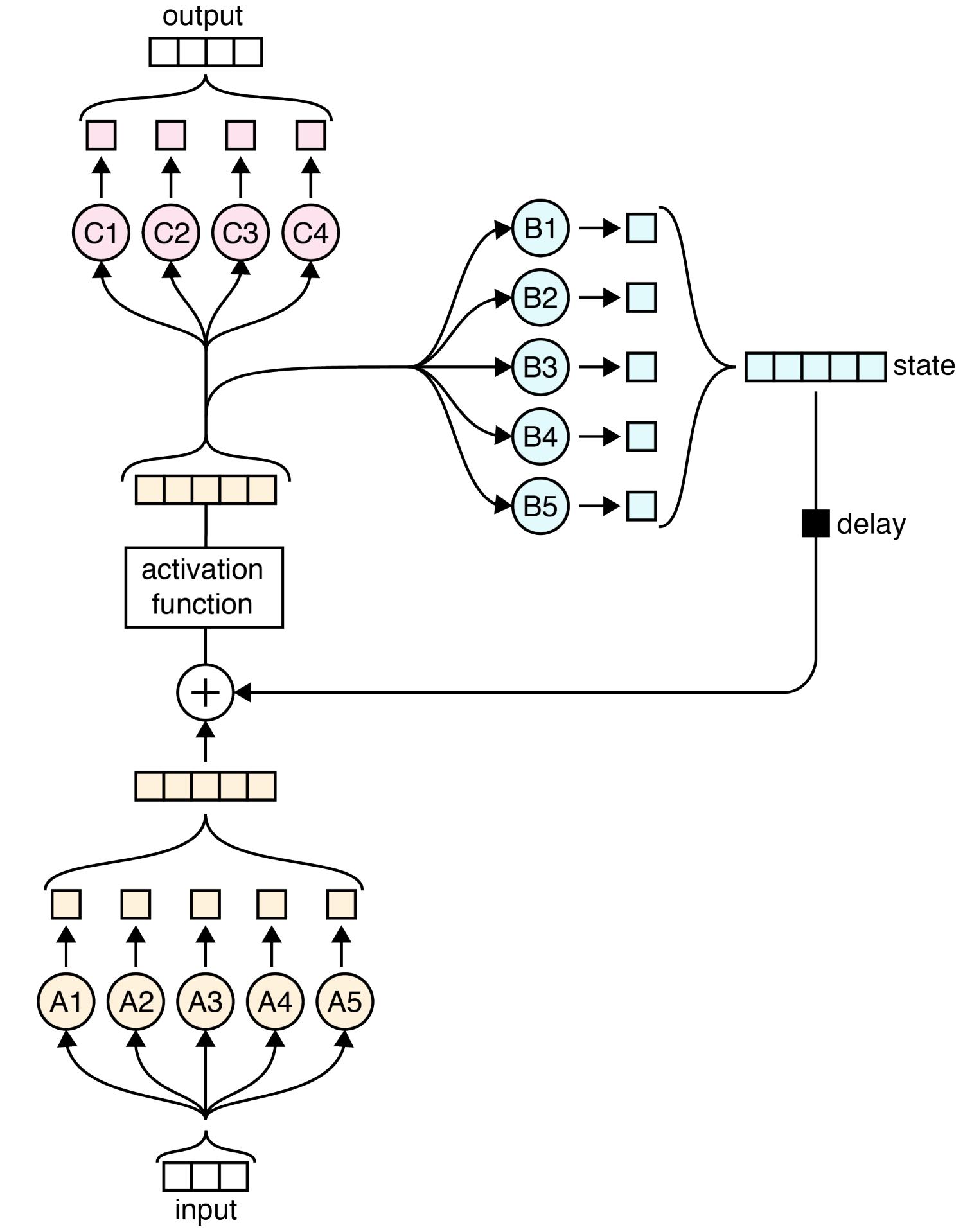
# Using Embeddings

```
model = Sequential()
model.add(Embedding(config.vocab_size,
                    config.embedding_dims,
                    input_length=config maxlen))
model.add(Conv1D(config.filters,
                 config.kernel_size,
                 padding='valid',
                 activation='relu'))
model.add(Flatten())
model.add(Dense(config.hidden_dims, activation='relu'))
model.add(Dense(1, activation='sigmoid'))

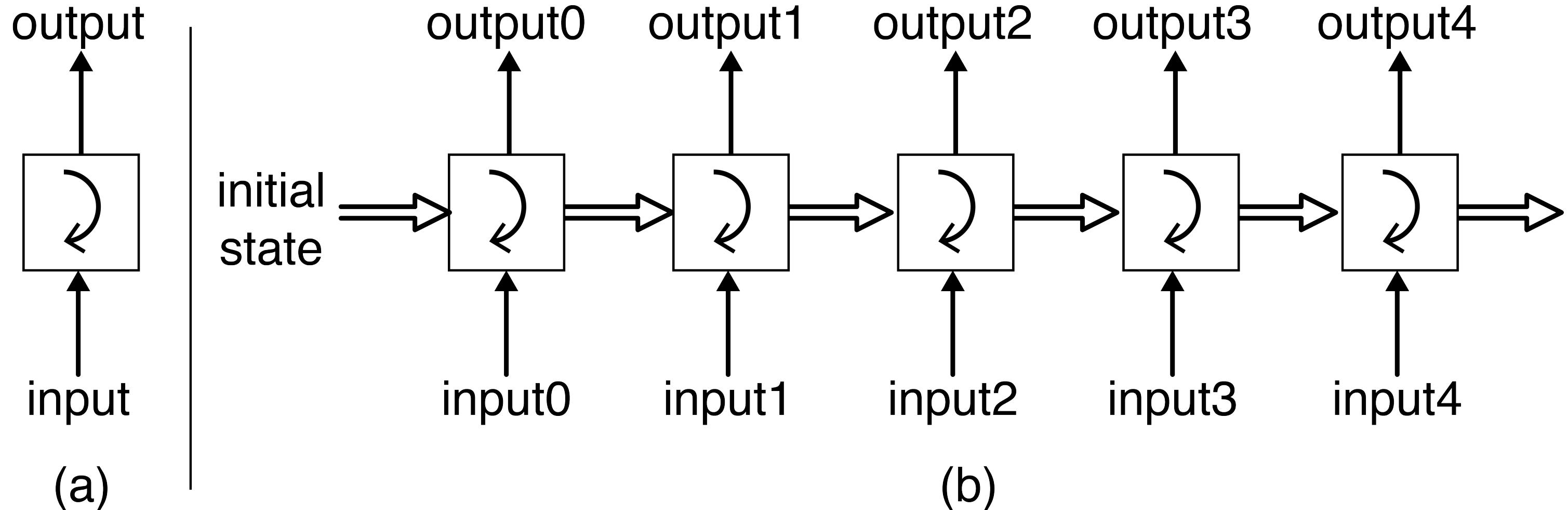
#OR

model.add(Embedding(config.vocab_size, 100,
input_length=config maxlen, weights=[embedding_matrix], trainable=False))
model.add(LSTM(config.hidden_dims, activation="sigmoid"))
```

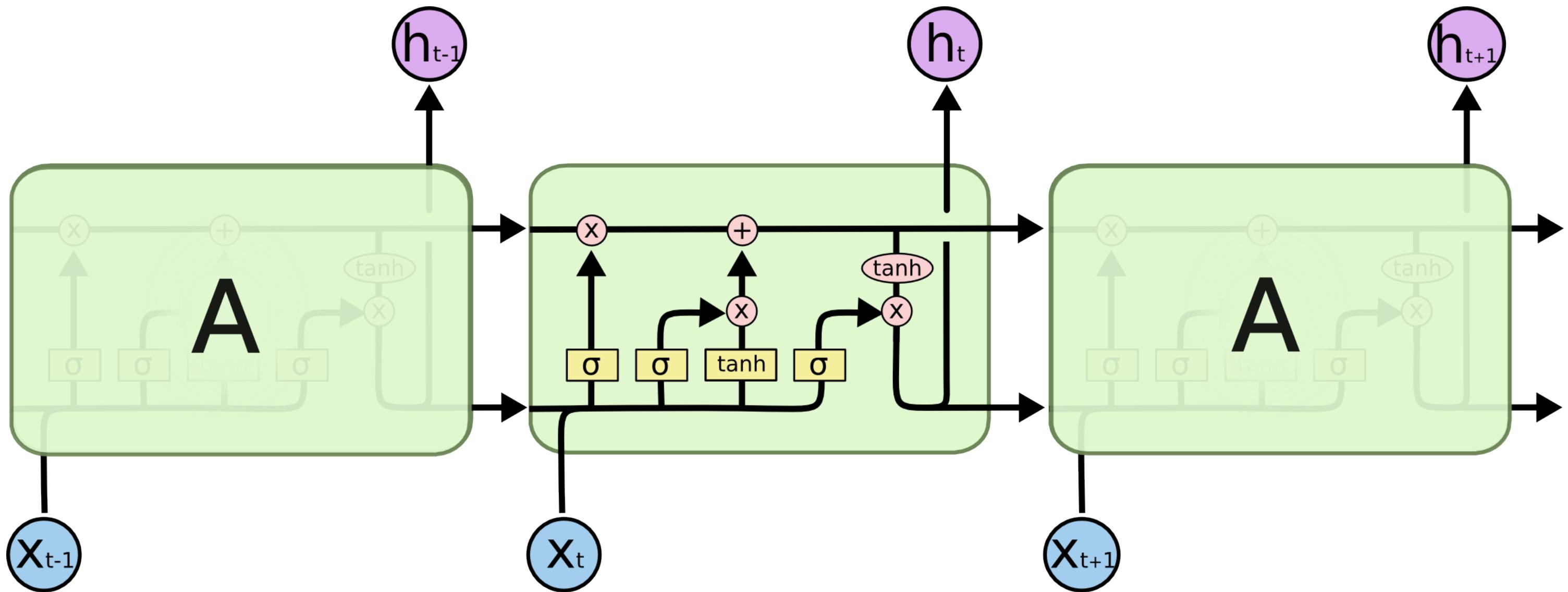
# Recurrent Neural networks

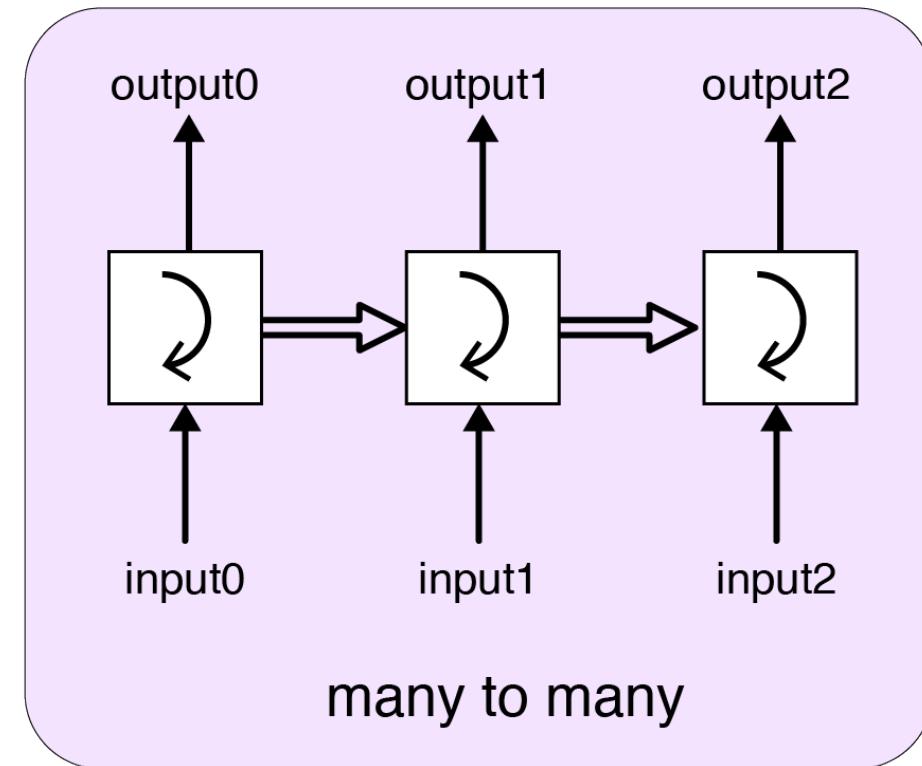
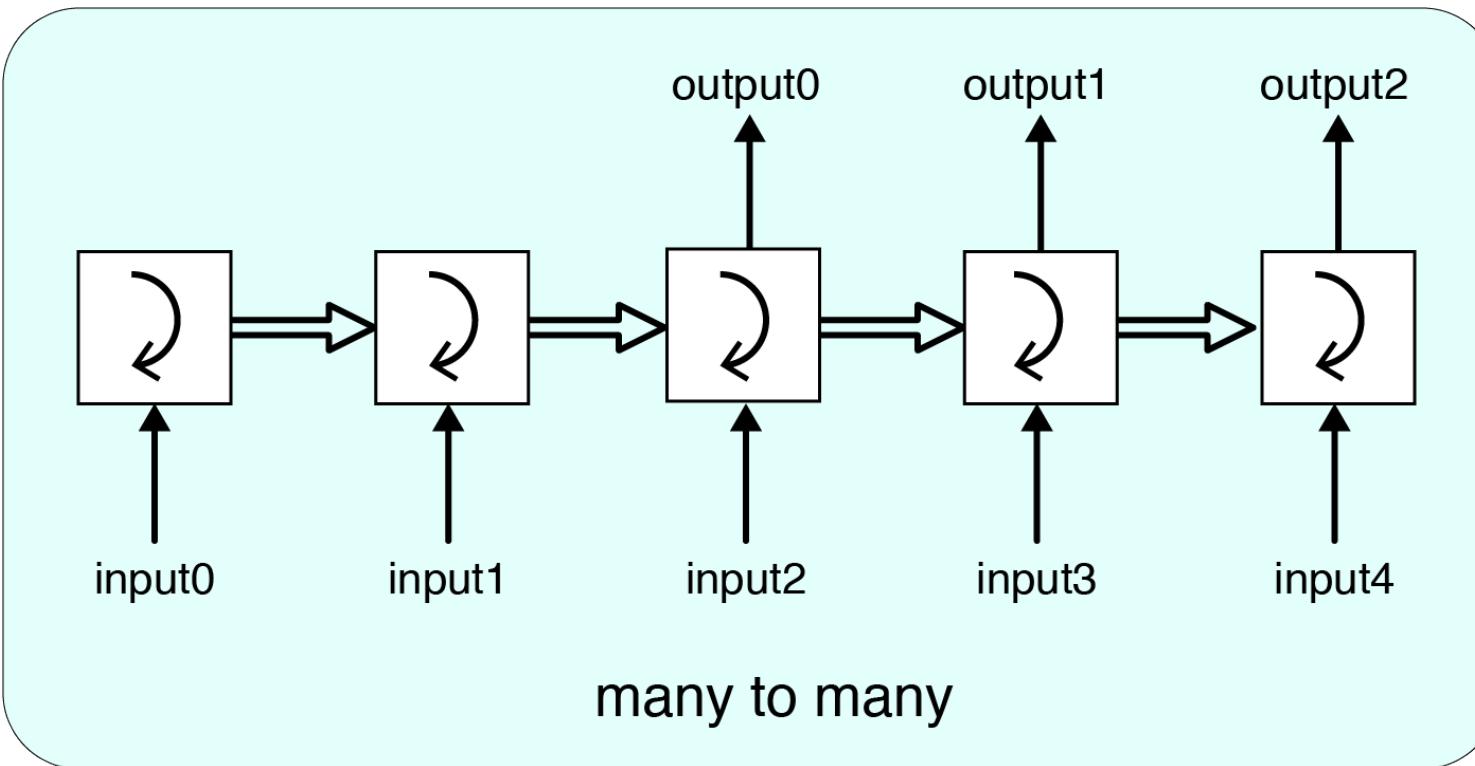
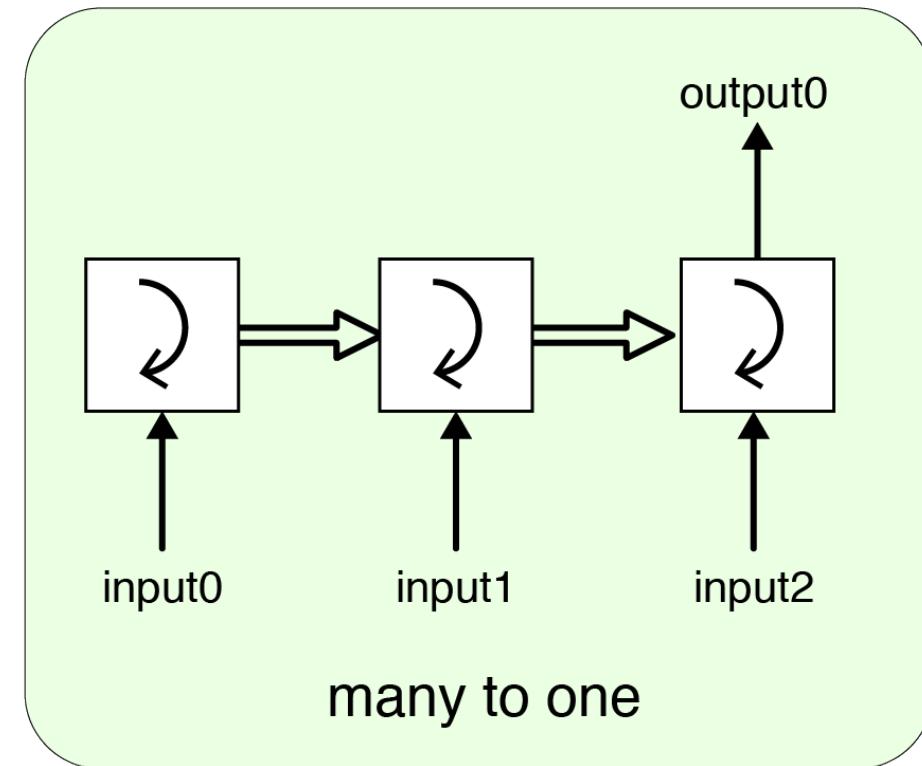
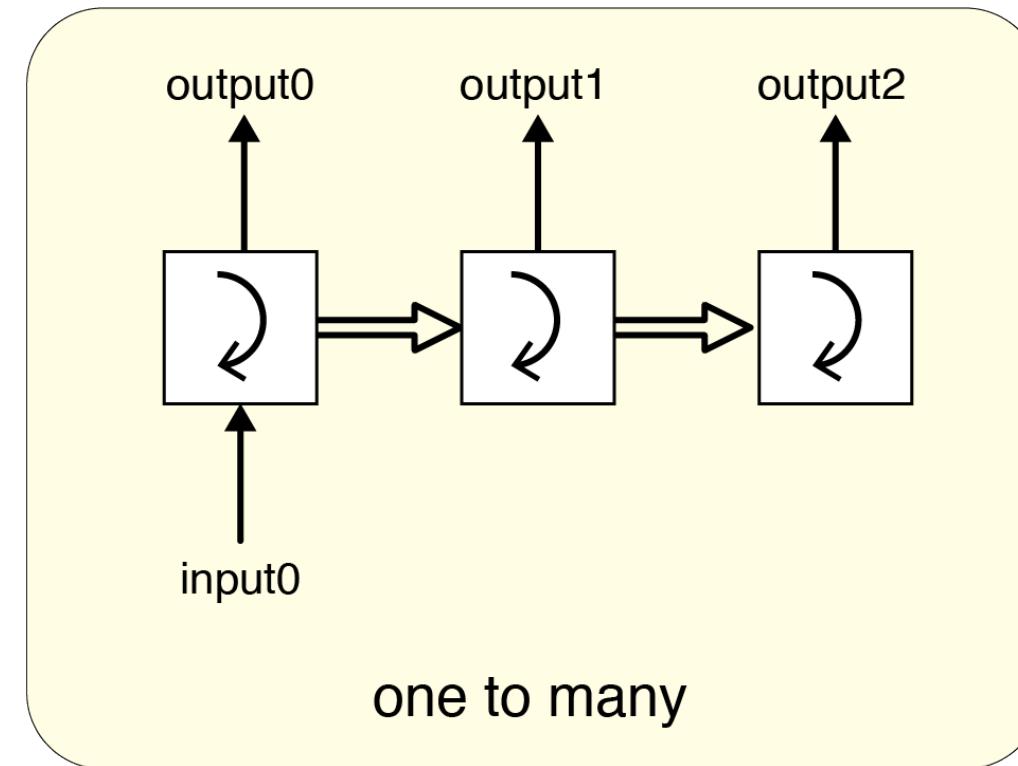
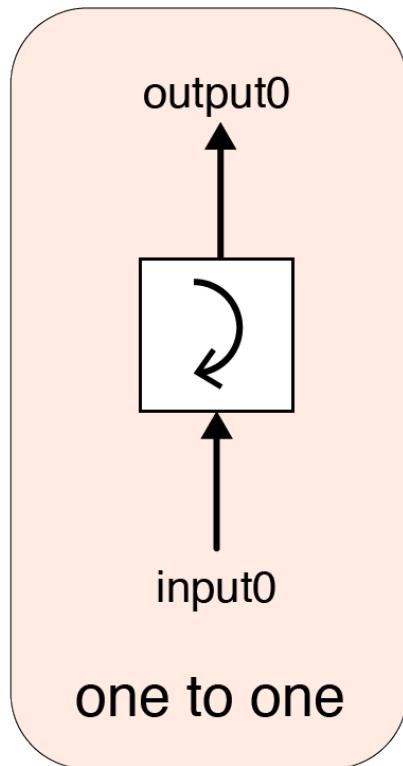


Unrolled...

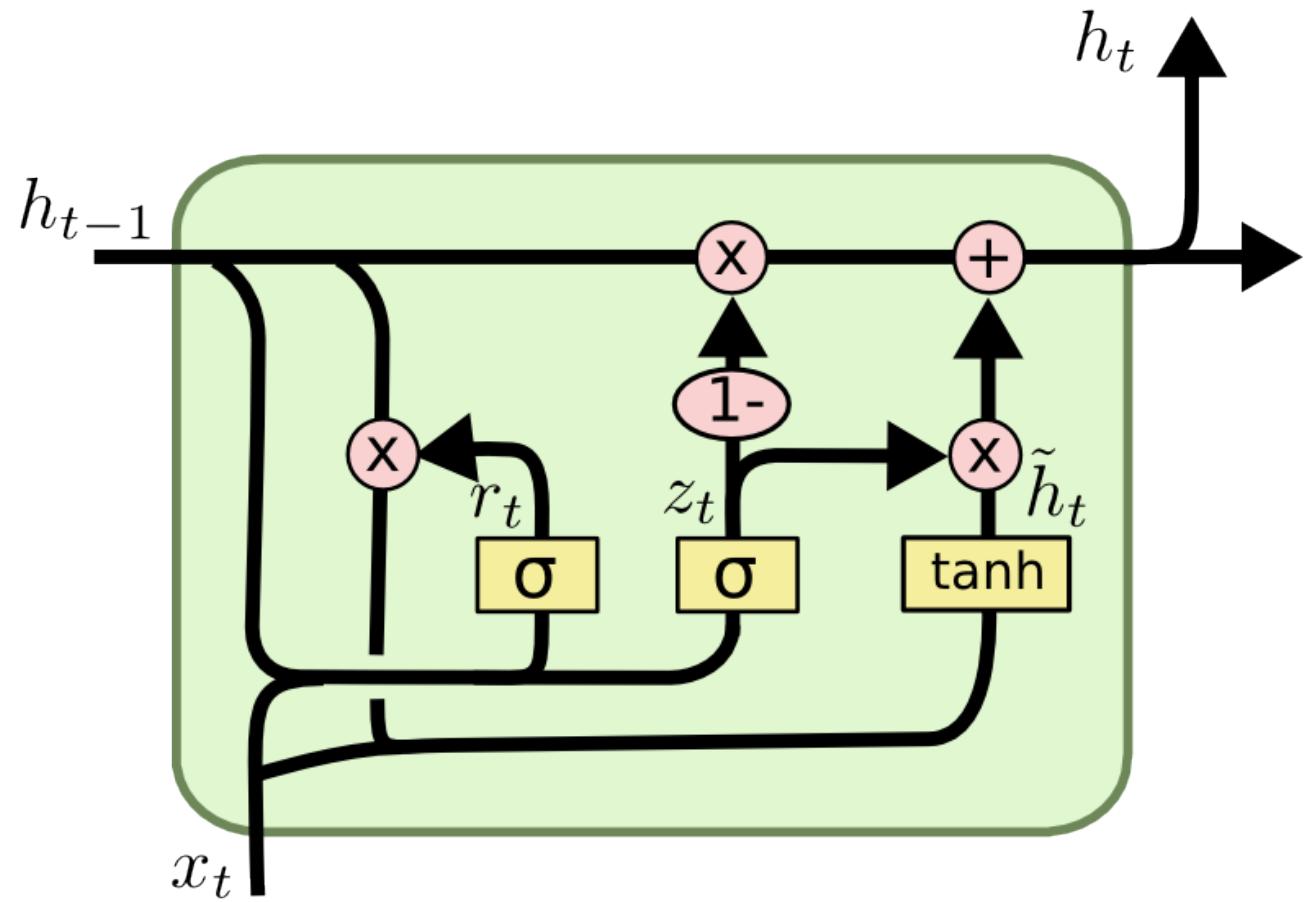


# LSTM for long term memory (vanishing gradients)





# GRU



$$z_t = \sigma (W_z \cdot [h_{t-1}, x_t])$$

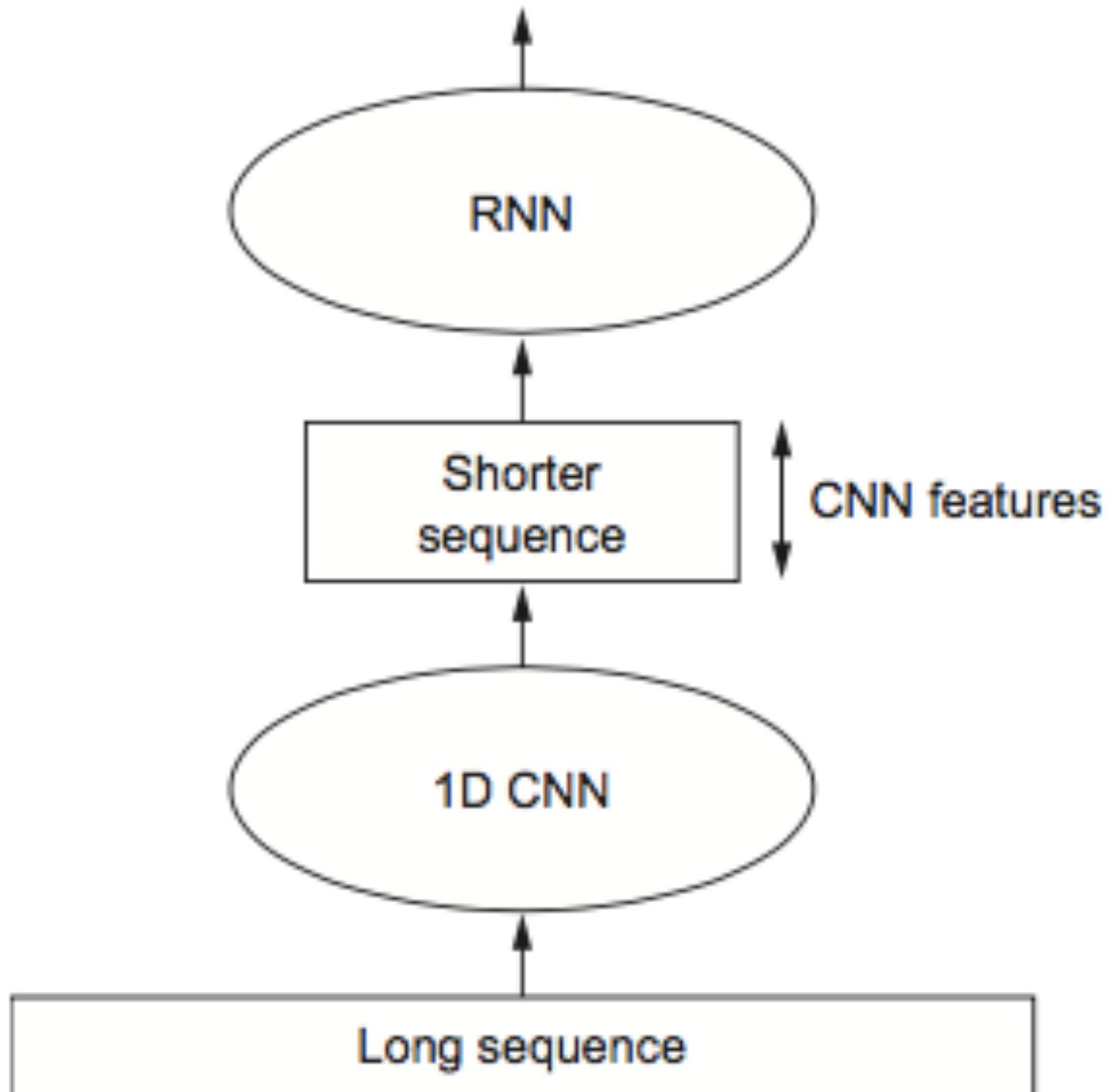
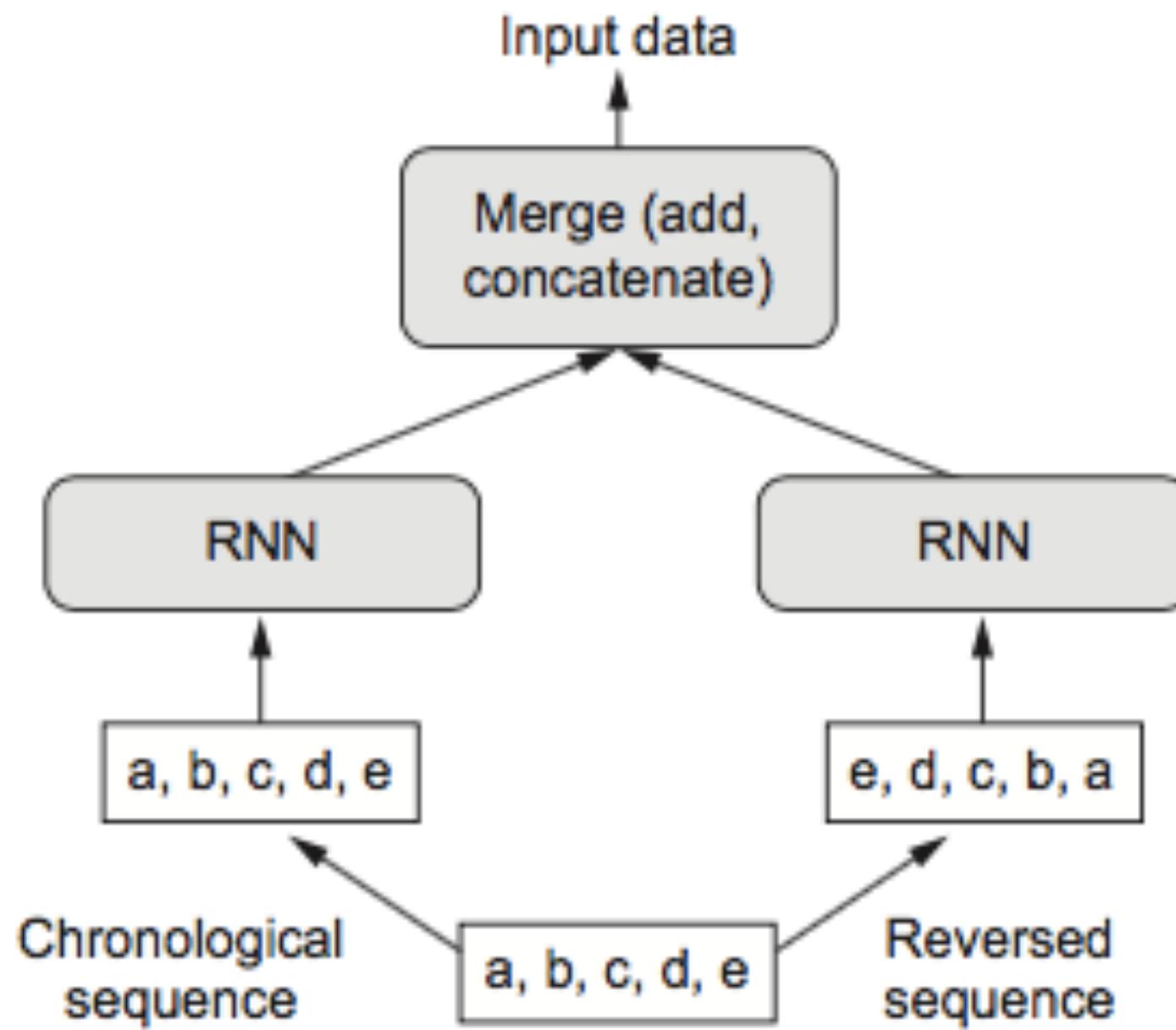
$$r_t = \sigma (W_r \cdot [h_{t-1}, x_t])$$

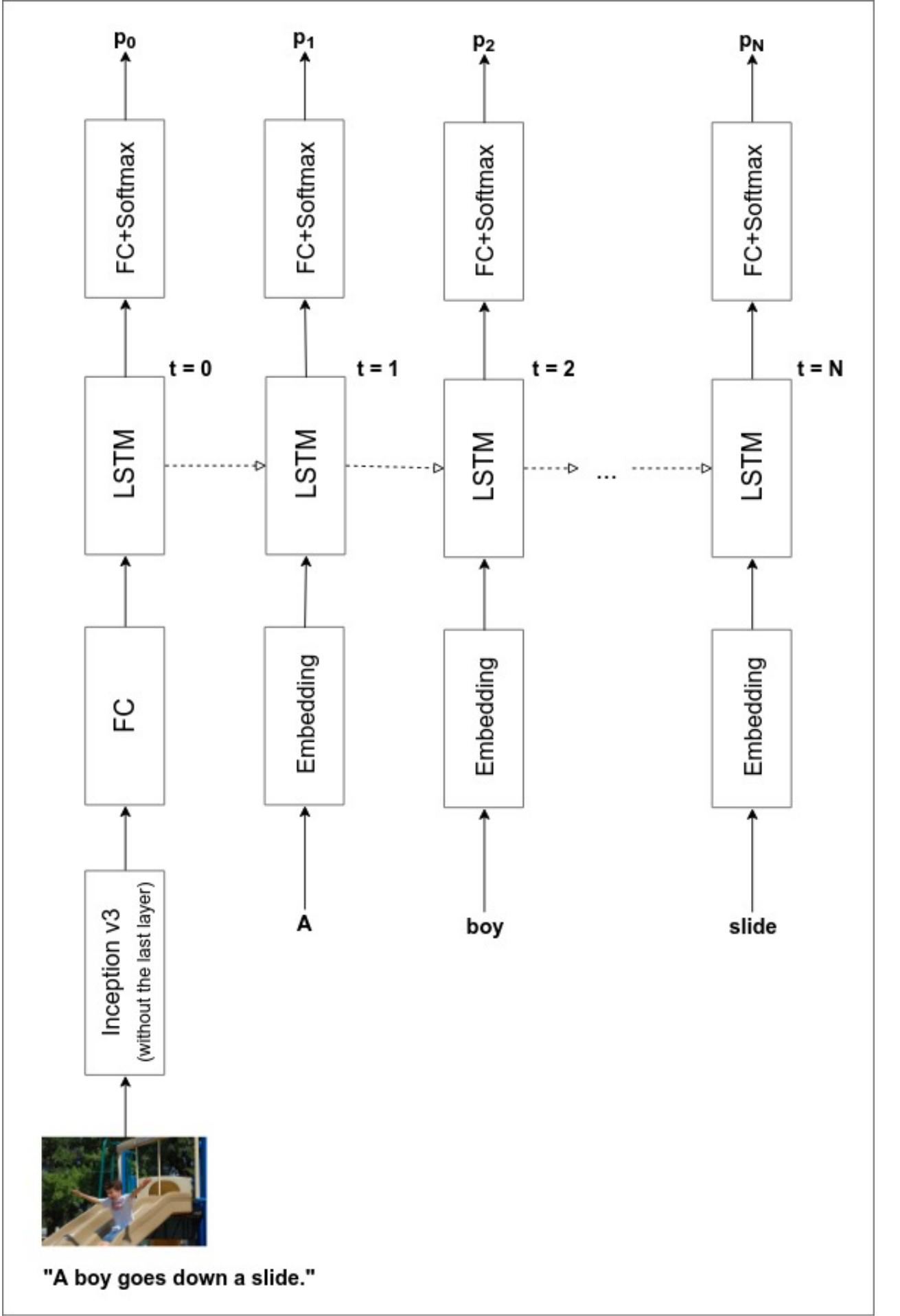
$$\tilde{h}_t = \tanh (W \cdot [r_t * h_{t-1}, x_t])$$

$$h_t = (1 - z_t) * h_{t-1} + z_t * \tilde{h}_t$$

# Other architectures

- cnn on embedding
- cnn-lstm
- stacked deep lstm
- bi-directional lstm
- cnn feeding into part of lstm (captioning)





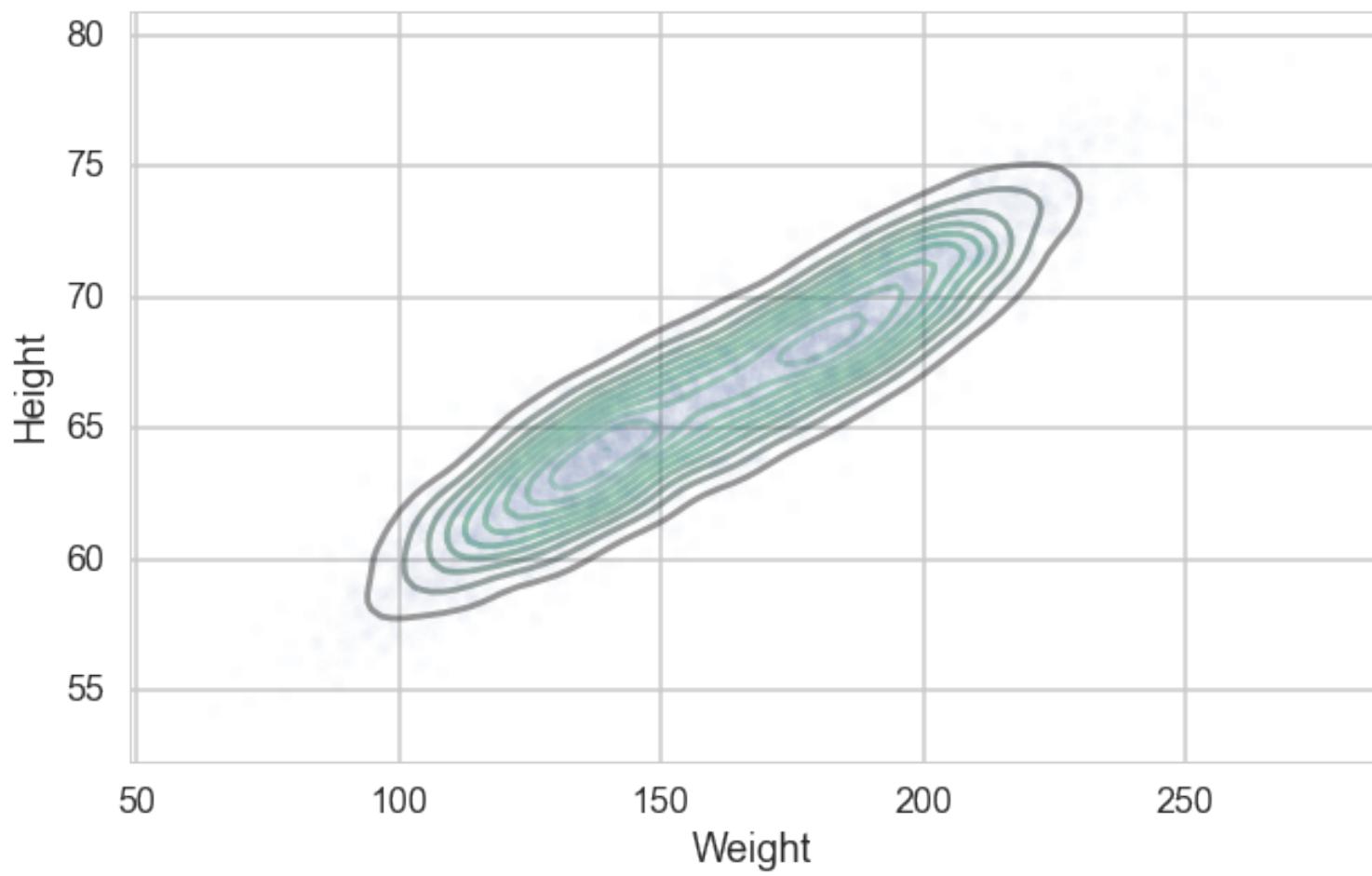
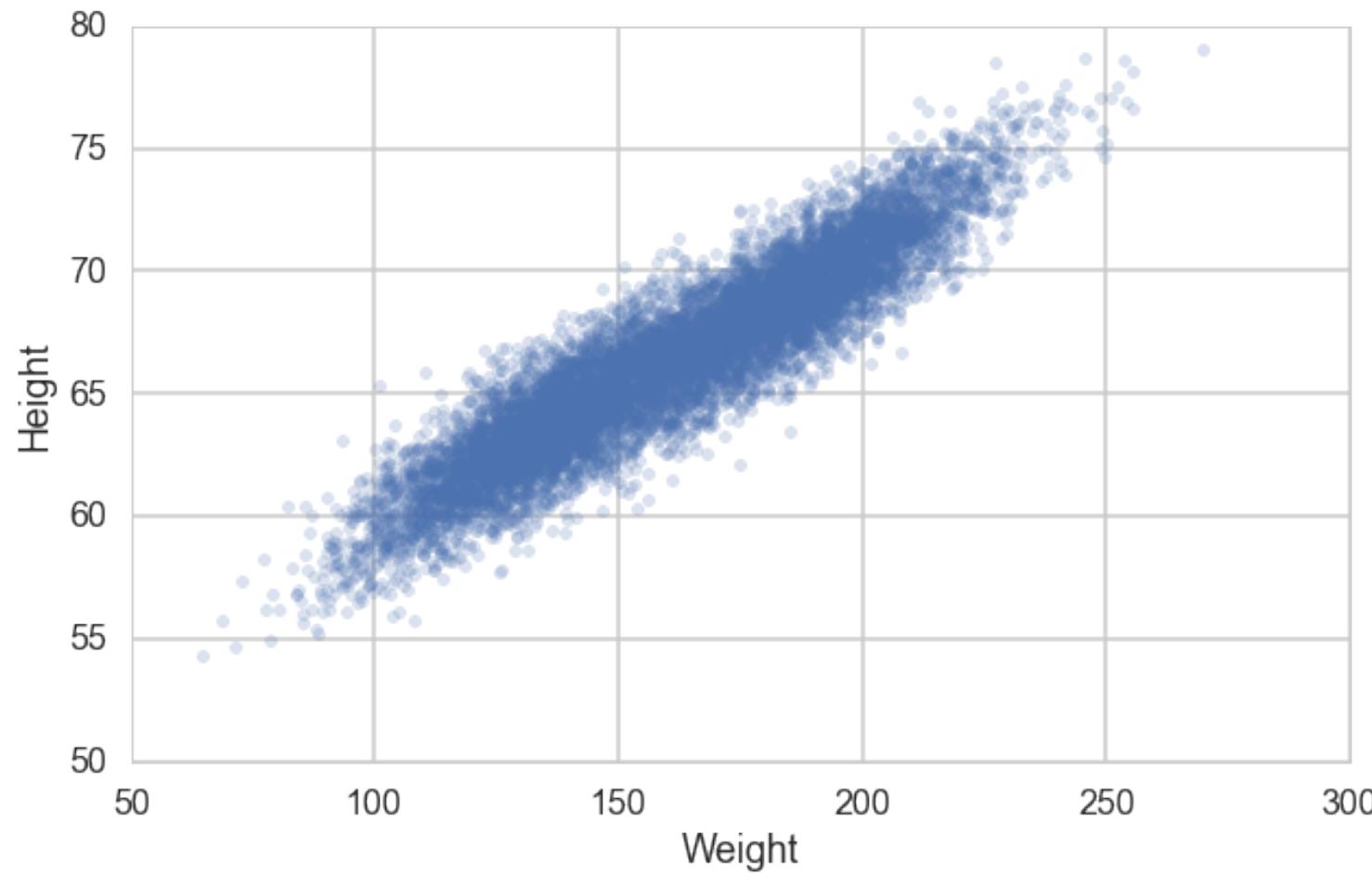
# Captioning

What other architectures could we use?

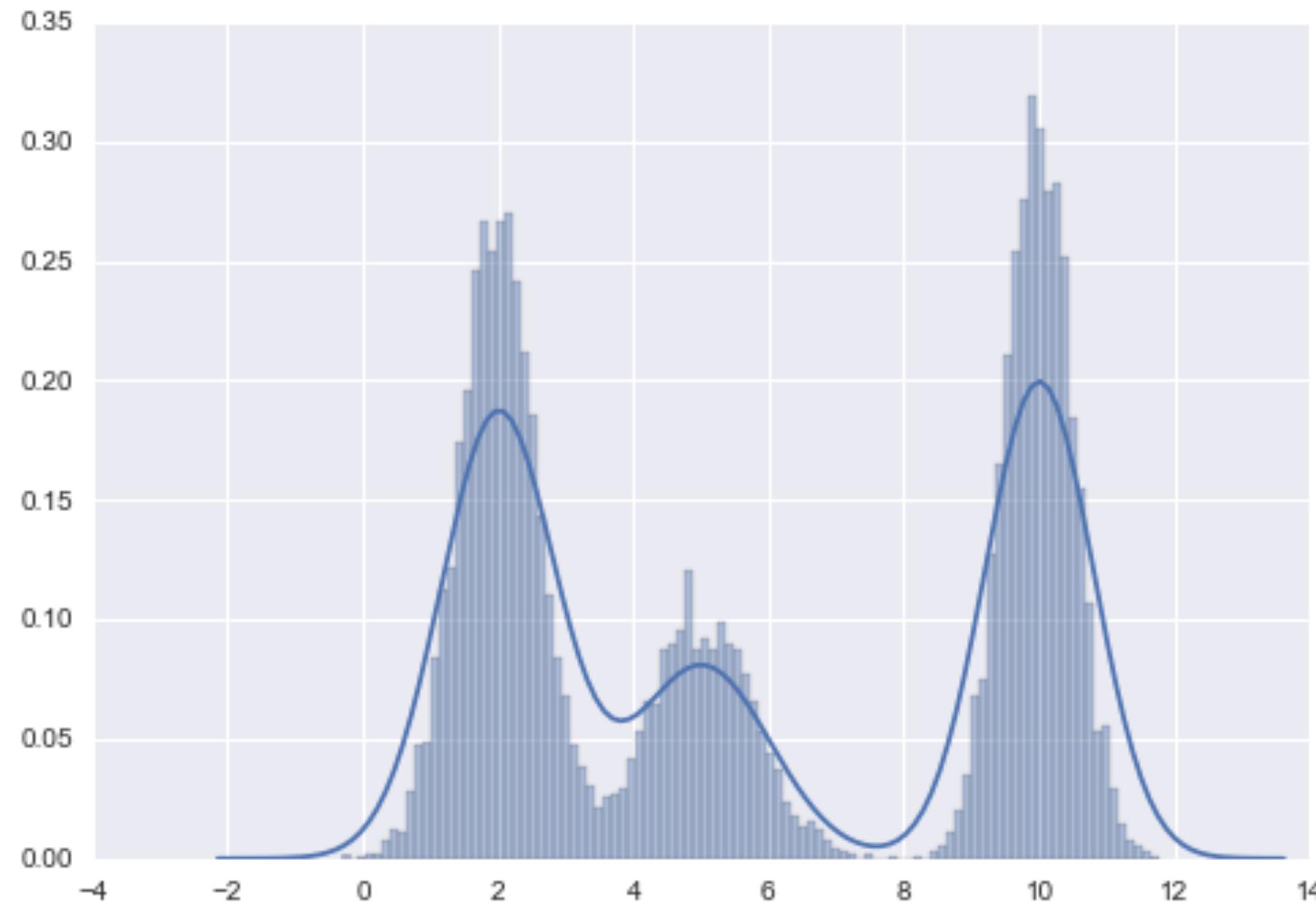
"A boy goes down a slide."

Generative  
Modeling

$p(x)$



# Big Question: Are these classes?



# Concrete Formulation of unsupervised learning

$$\begin{aligned} l(x|\lambda, \mu, \Sigma) &= \sum_{i=1}^m \log p(x_i|\lambda, \mu, \Sigma) \\ &= \sum_{i=1}^m \log \sum_z p(x_i|z_i, \mu, \Sigma) p(z_i|\lambda) \end{aligned}$$

Not Solvable analytically!

# Supervised vs Unsupervised Learning

In **Supervised Learning**, Latent Variables  $\mathbf{z}$  are observed.

In other words, we can write the full-data likelihood  $p(\mathbf{x}, \mathbf{z})$

In **Unsupervised Learning**, Latent Variables  $\mathbf{z}$  are hidden.

We can only write the observed data likelihood:

$$p(\mathbf{x}) = \sum_z p(\mathbf{x}, \mathbf{z}) = \sum_z p(\mathbf{z})p(\mathbf{x}|\mathbf{z})$$

COMBINE: Semi supervised learning.



In general

Representations are not classes

they are tangled, hierarchical complex things

but learning them

makes all the difference...

# From unsupervised learning

## ■ "Pure" Reinforcement Learning (cherry)

- ▶ The machine predicts a scalar reward given once in a while.
- ▶ **A few bits for some samples**

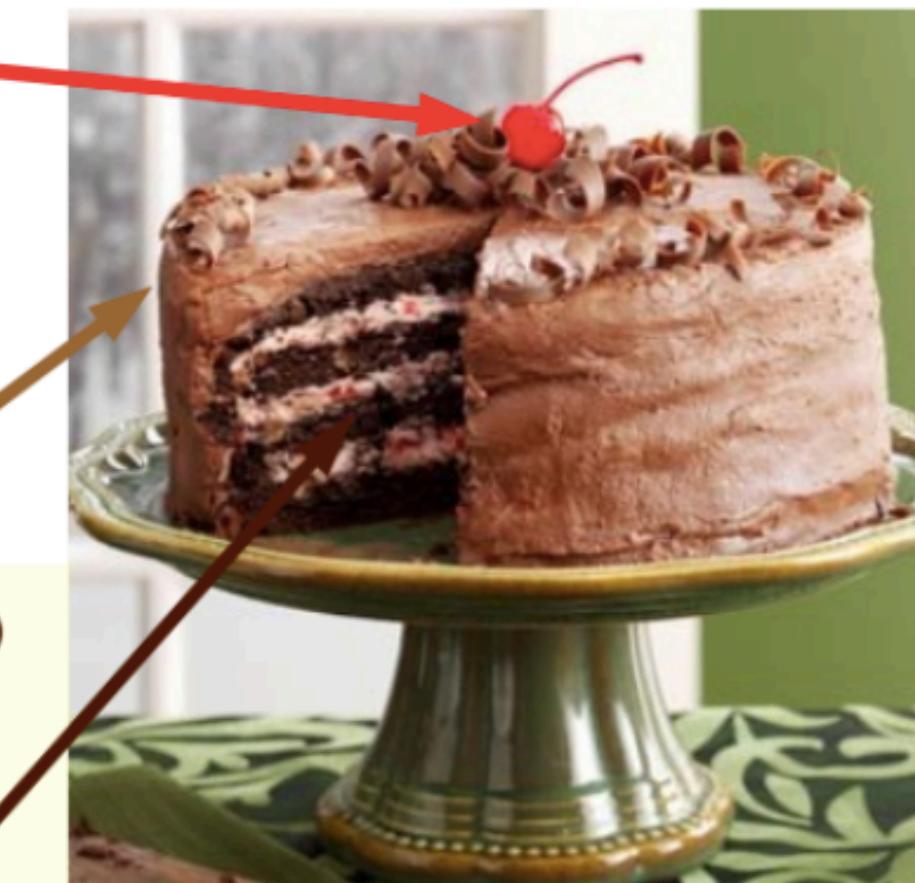
## ■ Supervised Learning (icing)

- ▶ The machine predicts a category or a few numbers for each input
- ▶ Predicting human-supplied data
- ▶ **10→10,000 bits per sample**

## ■ Unsupervised/Predictive Learning (cake)

- ▶ The machine predicts any part of its input for any observed part.
- ▶ Predicts future frames in videos
- ▶ **Millions of bits per sample**

■ (Yes, I know, this picture is slightly offensive to RL folks. But I'll make it up)



# To self-supervised learning

Y. LeCun

## How Much Information is the Machine Given during Learning?

- ▶ “Pure” Reinforcement Learning (**cherry**)
  - ▶ The machine predicts a scalar reward given once in a while.
  - ▶ **A few bits for some samples**
- ▶ Supervised Learning (**icing**)
  - ▶ The machine predicts a category or a few numbers for each input
  - ▶ Predicting human-supplied data
  - ▶ **10→10,000 bits per sample**
- ▶ Self-Supervised Learning (**cake génoise**)
  - ▶ The machine predicts any part of its input for any observed part.
  - ▶ Predicts future frames in videos
  - ▶ **Millions of bits per sample**

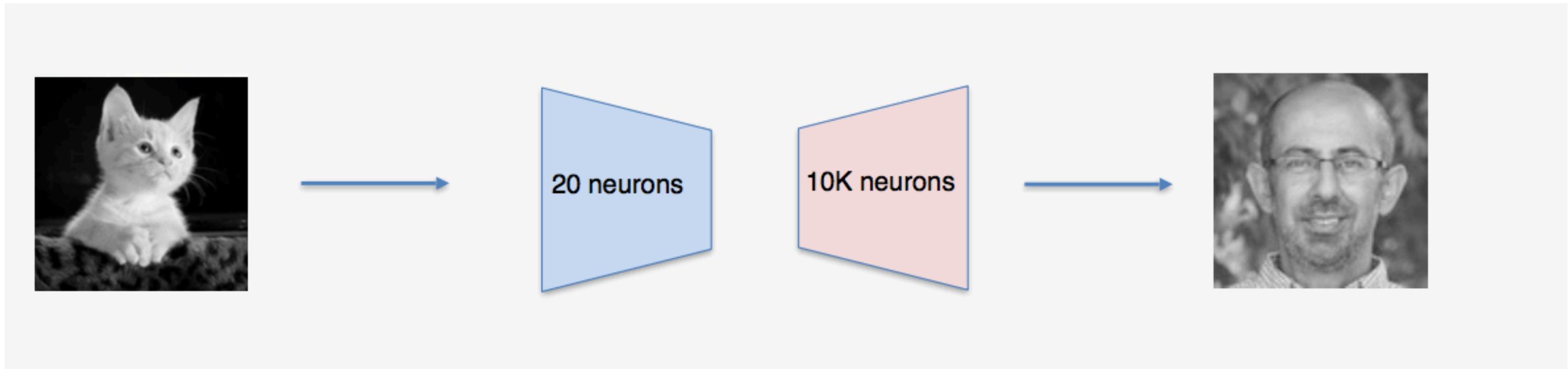


© 2019 IEEE International Solid-State Circuits Conference

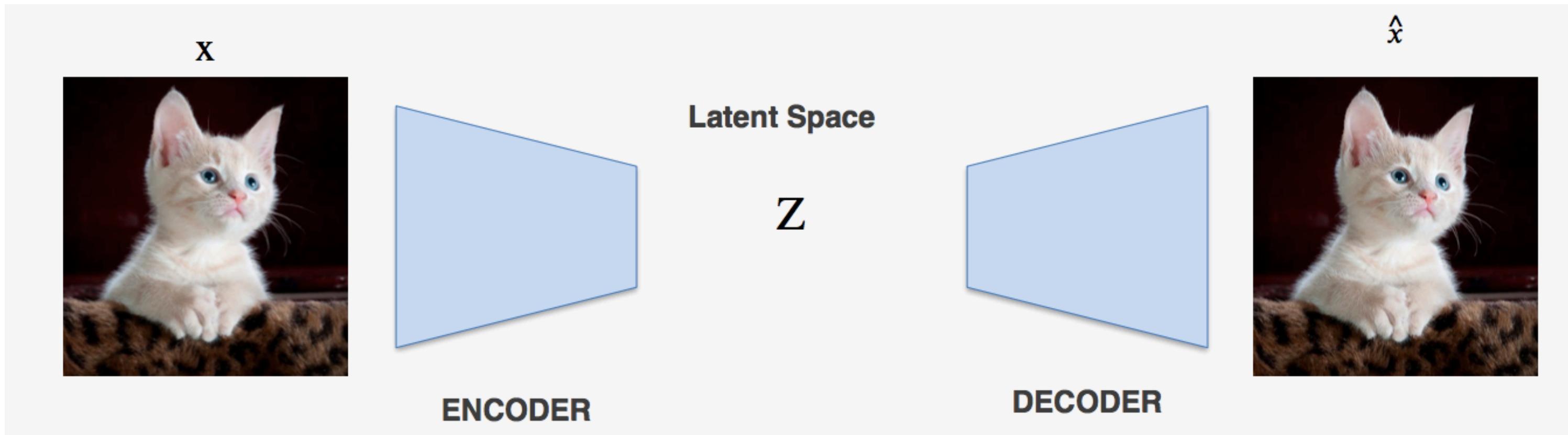
1.1: Deep Learning Hardware: Past, Present, & Future

59

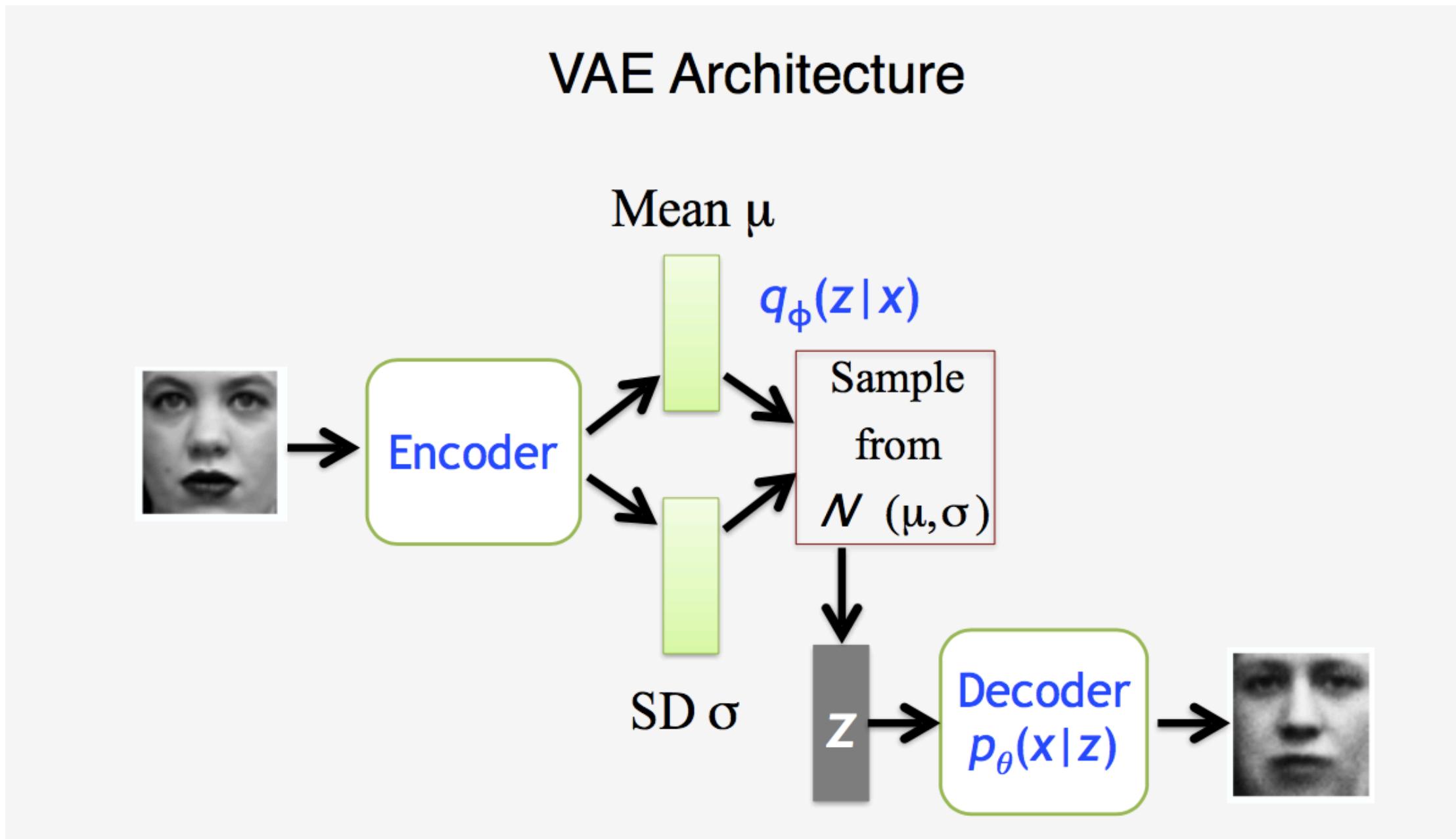
# General Encoder-Decoder Architecture



# Autoencoder



# Variational Autoencoder



# VAE Loss

Reconstruction Loss

Proposal distribution  
should resemble a Gaussian

$$-\mathbb{E}_{z \sim q_{\phi}(z|x)} \log(p_{\theta}(x|z)) + KL(q_{\phi}(z|x) \parallel p_{\theta}(z))$$

# Moving from one rep to the other

# VAE: Deep Generative Models

- simply not possible to do inference in large models
- inference in neural networks: understanding robustness, etc
- hierarchical neural networks
- Mixture density networks: mixture parameters are fitted using ANNs
- extension to generative semisupervised learning
- variational autoencoders

Where do we go from here?

What should you do?

What should you read? See?

The advanced course

# What should you do?

- There is no substitute for coding.
- We are there to help. Please log onto <https://discourse.univ.ai>, so we can have conversations (use your github id).
- In this next week, choose one of the hackings for the course and work (more) on it.
- we will have a hack every 2 weeks (so 2 a month) until the next basics (next 3 months). we'll discuss on Discourse
- suggest topics if interested!

## Come TA for us

- we are looking for TAs for the next basics (hybrid on-site/online, october/november) and the advanced (online, november/december). Ping us if you want to do this!
- TAing is the best way to take your learning to a new level, nothing forces you to understand something like having to explain it.
- you will have the opportunity to develop material, and this will help your understanding
- TA training will be provided

# Binge Watching/Reading (watch at 1.5x)

- The first session of fast.ai
- Pattern Recognition and Machine Learning, Bishop
- Deep Learning, Eugene Charniak
- Deep Learning, Andrew Glassner
- Andrew Ng's course is an oldie but goodie. Machine Learning Yearning PDF Document
- Elements of Statistical Learning

# The advanced course

- more applications: super-resolution, image segmentation. etc Unets.
- deep unsupervised learning: GANs, more on autoencoders, autoregressive models, flow models
- Transfer learning in NLP with Bert, Elmo, and ULM-Fit
- model deployment and experimentation
- seq2seq, attention
- bayesian statistics

**Get on Discourse!**