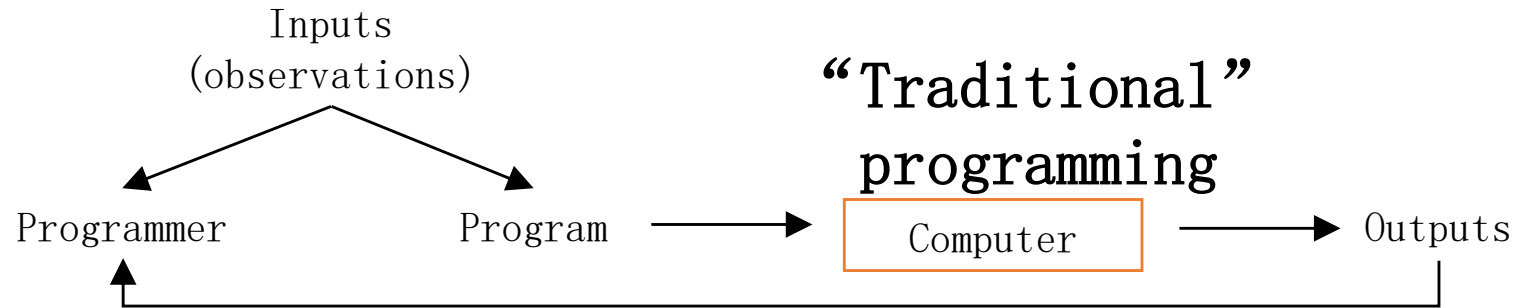# 机器学习和量化交易实战

第四讲

# Outline

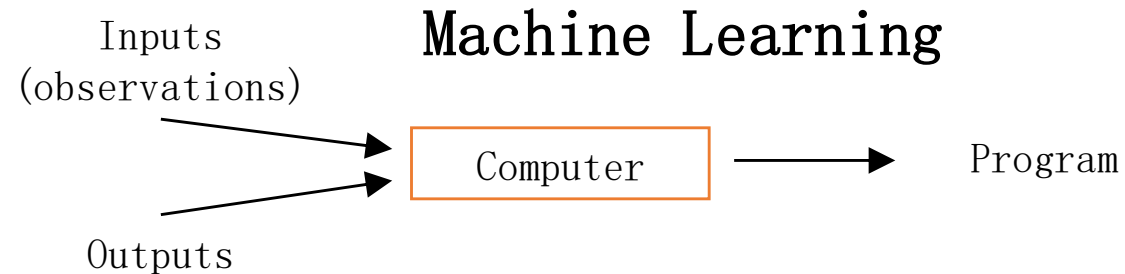- From OLS to kernel machines and beyond
  - OLS
  - Ridge
  - L a s s o
  - Kernels
  - Cross-validation
  - Hands on: sklearn

# What is Machine Learning?

Inputs
(observations)

"Traditional"
programming

Programmer        Program        Computer        Outputs

*Machine Learning is the field of study that gives computers the ability to learn without being explicitly programmed.*
*-- Arthur Samuel (1959)*

Inputs
(observations)

Machine Learning

Computer        Program

Outputs

# Examples of Machine Learning
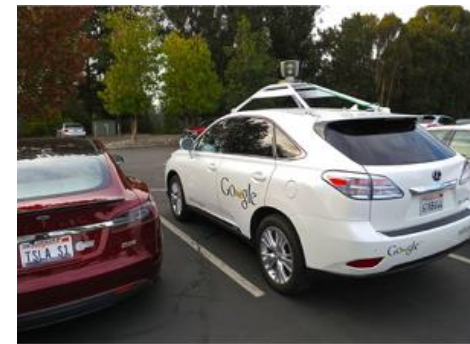
https://flic.kr/p/5BLW6G [CC BY 2.0]

http://commons.wikimedia.org/wiki/File:American_book_company_1916._letter_envelope-2.JPG#filelinks [public domain]

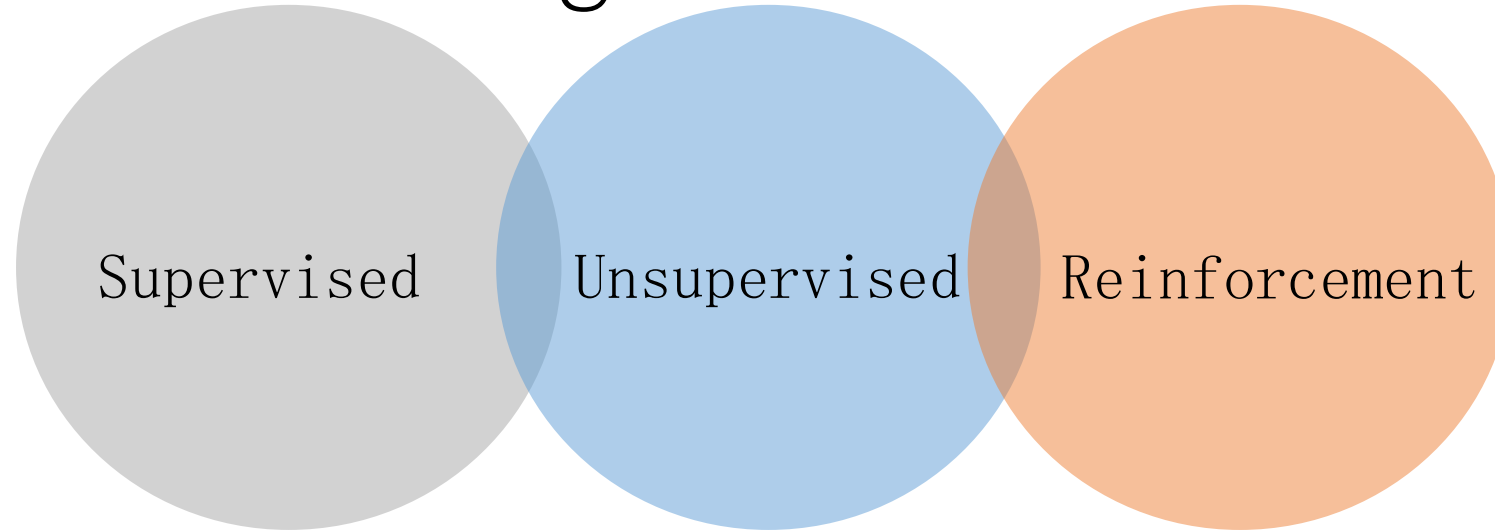http://commons.wikimedia.org/wiki/File:Netflix_logo.svg [public domain]

And many, many more …

By Steve Jurvetson [CC BY 2.0]
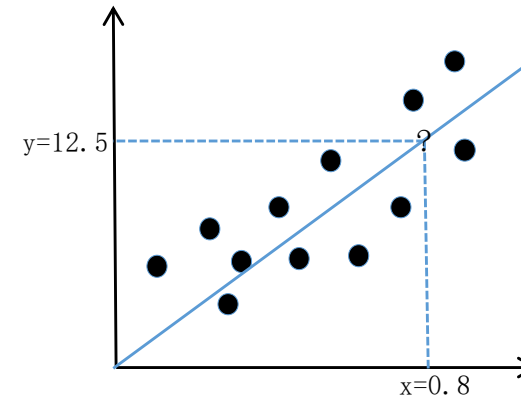
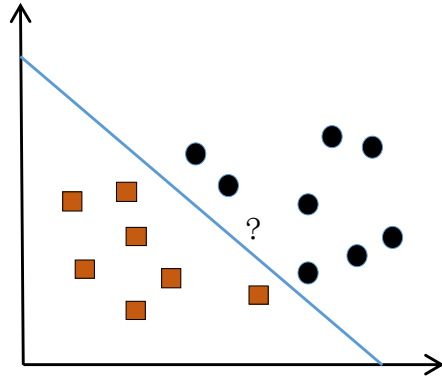# 3 Types of Learning

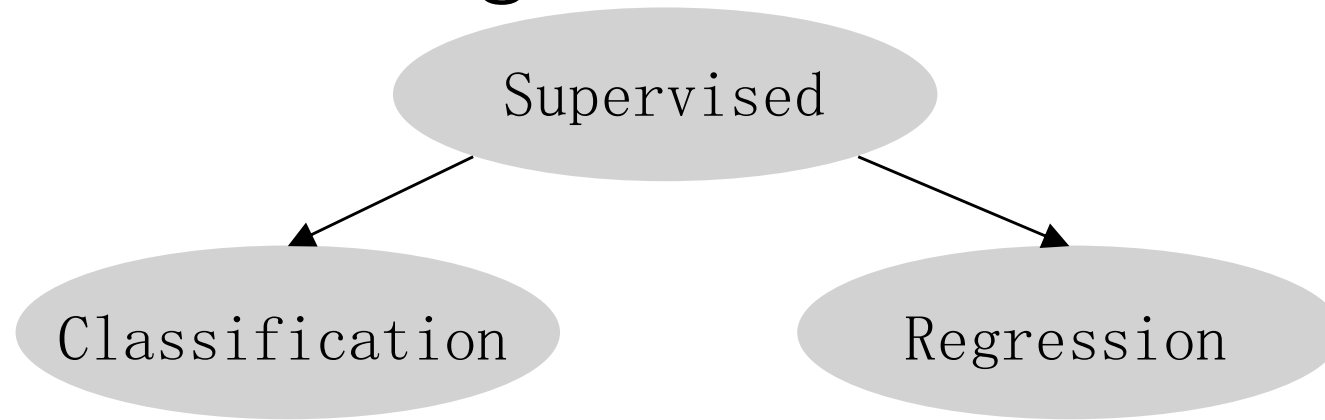Supervised    Unsupervised    Reinforcement

- Learning from labeled data
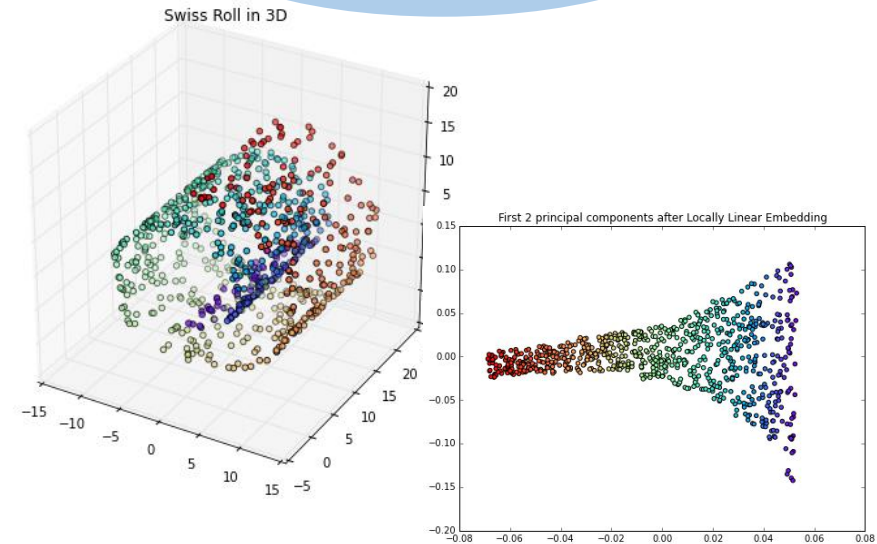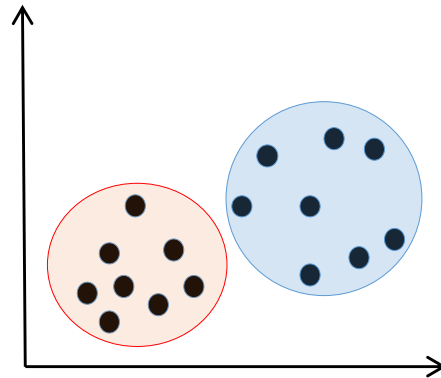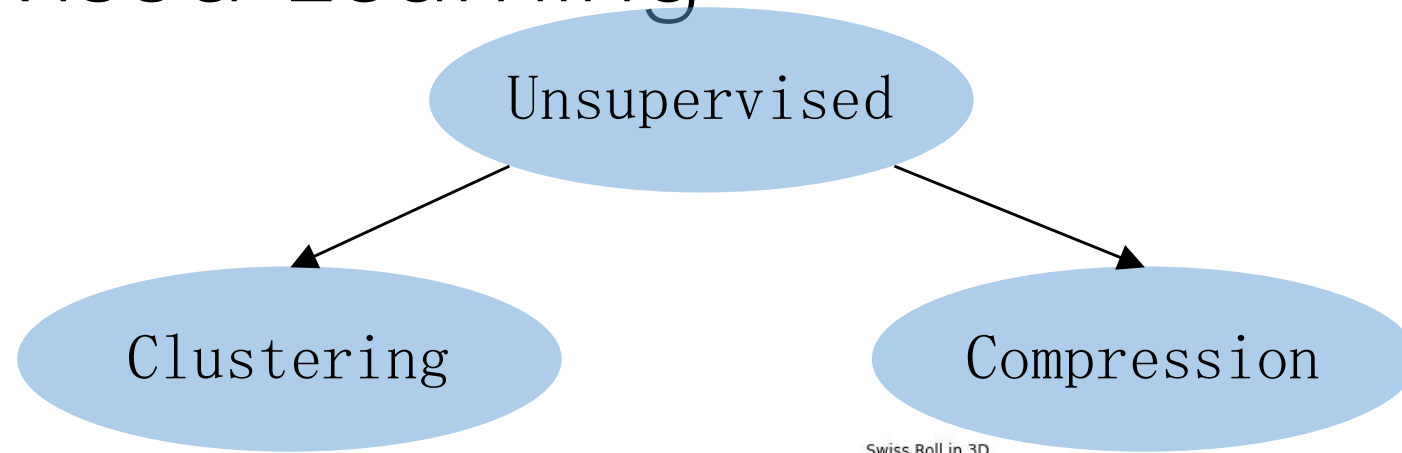- E.g., Spam classification

- Discover structure in unlabeled data
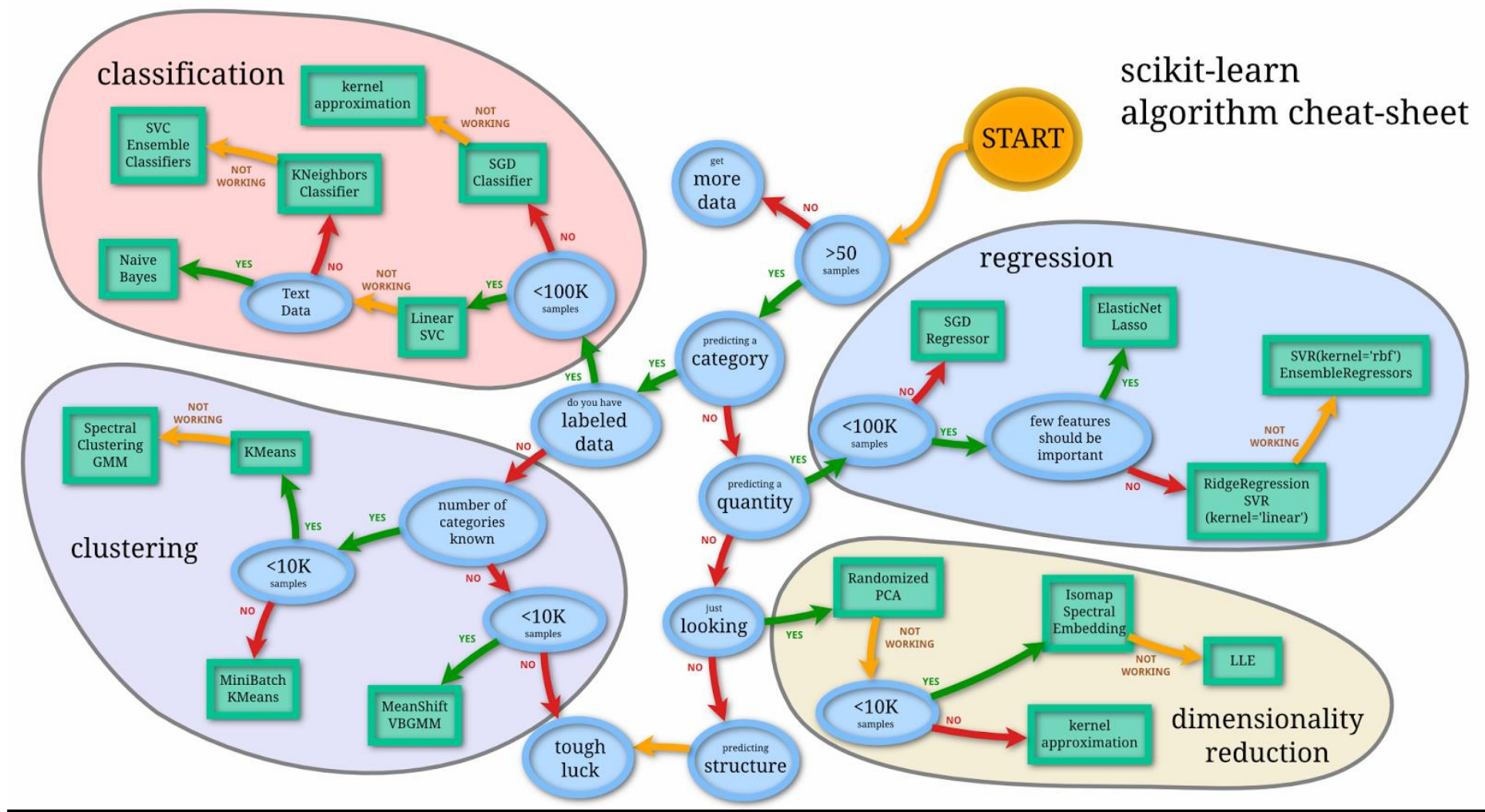- E.g., Document clustering

- Learning by "doing" with delayed reward
- E.g., Chess computer

# Supervised Learning

# Unsupervised Learning

scikit-learn algorithm cheat-sheet

**classification**

- kernel approximation
- SVC Ensemble Classifiers
- KNeighbors Classifier
- SGD Classifier
- Naive Bayes
- Text Data
- Linear SVC
- <100K samples

**START**

- get more data
- >50 samples
- predicting a category
- do you have labeled data

**regression**

- SGD Regressor
- ElasticNet Lasso
- SVR(kernel='rbf') EnsembleRegressors
- <100K samples
- few features should be important
- RidgeRegression SVR (kernel='linear')
- predicting a quantity

**clustering**

- Spectral Clustering GMM
- KMeans
- number of categories known
- <10K samples
- MiniBatch KMeans
- MeanShift VBGMM
- <10K samples

**dimensionality reduction**

- Randomized PCA
- Isomap Spectral Embedding
- LLE
- <10K samples
- kernel approximation
- just looking
- predicting structure
- tough luck

8

# The simplest Sklearn workflow

```
train_x, train_y, test_x, test_y = getData()

model = somemodel()
model.fit(train_x,train_y)
predictions = model.predict(test_x)

score = score_function(test_y, predictions)
```
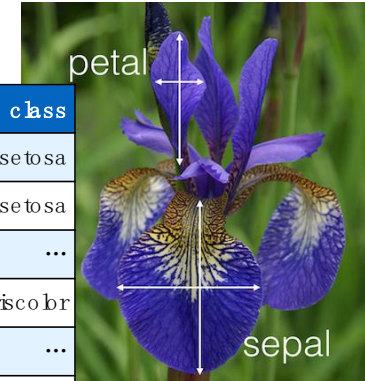
# Flower Classification

Iris-Setosa

Iris-Setosa

Iris-Versicolor

# Data Representation

https://archive.ics.uci.edu/ml/datasets/Iris

Instances (samples, observations)

| | sepal_length | sepal_width | petal_length | petal_width | class |
|---|---|---|---|---|---|
| 1 | 5.1 | 3.5 | 1.4 | 0.2 | setosa |
| 2 | 4.9 | 3.0 | 1.4 | 0.2 | setosa |
| ... | ... | ... | ... | ... | ... |
| 50 | 6.4 | 3.2 | 4.5 | 1.5 | veriscolor |
| ... | ... | ... | ... | ... | ... |
| 150 | 5.9 | 3.0 | 5.1 | 1.8 | virginica |

petal

sepal

Features (attributes, dimensions)

Classes (targets)

$$\mathbf{X} = \begin{bmatrix} x_{11} & x_{12} & \cdots & x_{1D} \\ x_{21} & x_{22} & \cdots & x_{2D} \\ x_{31} & x_{32} & \cdots & x_{3D} \\ \vdots & \vdots & \ddots & \vdots \\ \vdots & \vdots & \ddots & \vdots \\ x_{N1} & x_{N2} & \cdots & x_{ND} \end{bmatrix}$$

$$\mathbf{y} = [y_1, y_2, y_3, \cdots y_N]$$

11

```
In [2]: from sklearn.datasets import load_iris
        iris = load_iris()
```

The resulting dataset is a `Bunch` object: you can see what's available using the method `keys()`:

```
In [3]: iris.keys()

Out[3]: dict_keys(['target_names', 'data', 'feature_names', 'DESCR', 'target'])
```
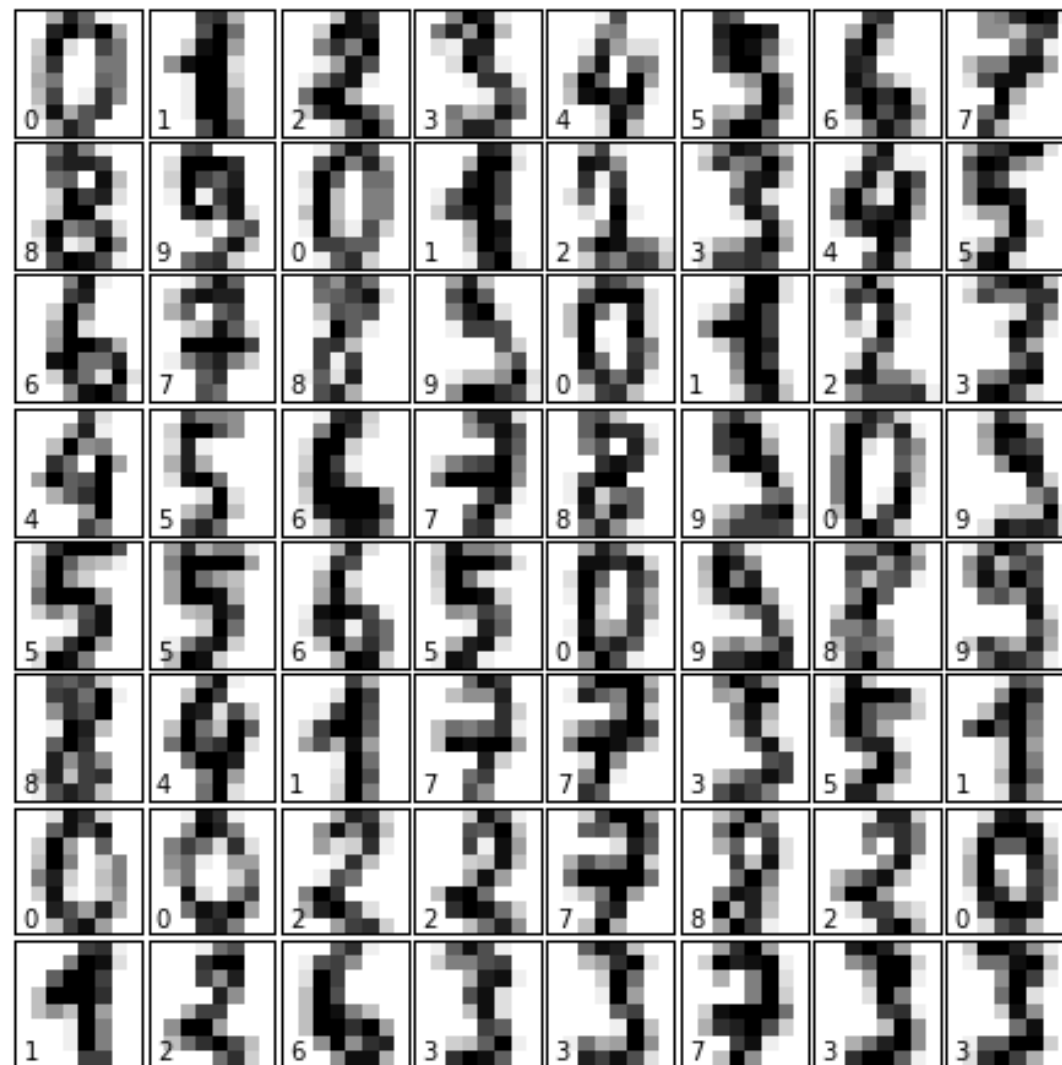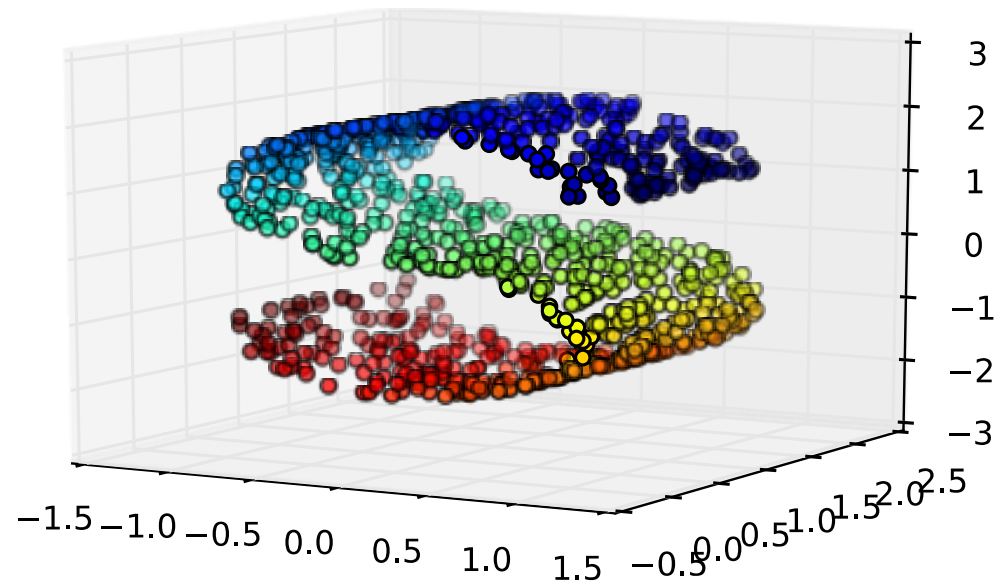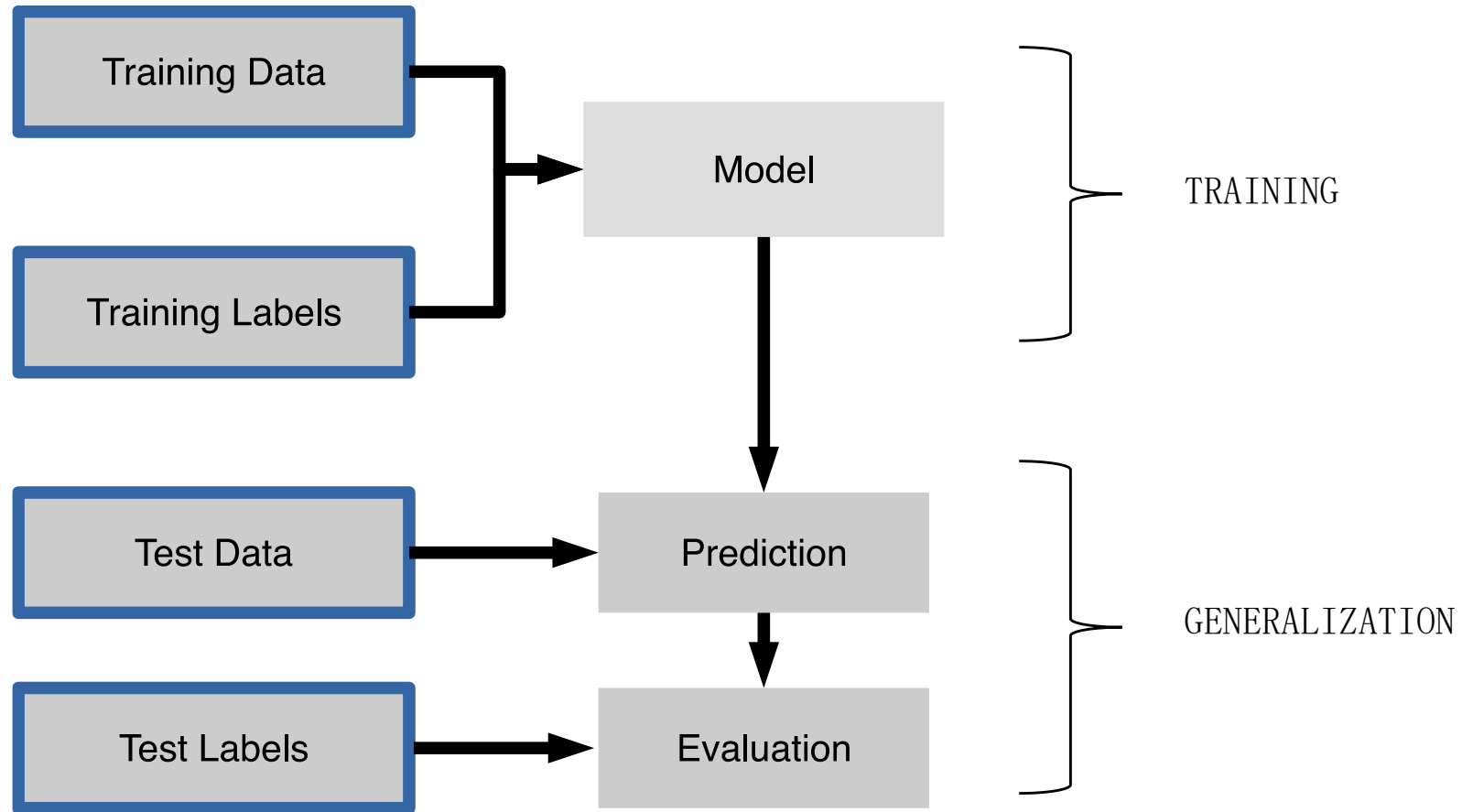
Iris-Setosa

Iris-Setosa

Iris-Versicolor

# Digits

# Generating Synthetic Data

```
from sklearn.datasets import make_...
```

# Supervised Workflow

# Supervised Workflow



estimator.fit(X_train, y_train)

TRAINING

estimator.predict(X_test)

GENERALIZATION

estimator.score(X_test, y_test)

# Regression Shrinkage and Selection via the Lasso

# Regularization

All the answers so far are of the form

$$\widehat{\boldsymbol{\theta}} = (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{y}$$

They require the inversion of $\mathbf{X}^T\mathbf{X}$. This can lead to problems if the system of equations is poorly conditioned. A solution is to add a small element to the diagonal:

$$\widehat{\boldsymbol{\theta}} = (\mathbf{X}^T\mathbf{X} + \delta^2 I_d)^{-1}\mathbf{X}^T\mathbf{y}$$

This is the ridge regression estimate. It is the solution to the following **regularised quadratic cost function**

$$J(\boldsymbol{\theta}) = (\mathbf{y} - \mathbf{X}\boldsymbol{\theta})^T(\mathbf{y} - \mathbf{X}\boldsymbol{\theta}) + \delta^2\boldsymbol{\theta}^T\boldsymbol{\theta}$$

# Derivation

$$\frac{\partial}{\partial \theta} J(\theta) = \frac{\partial}{\partial \theta} \left\{ (Y - X\theta)^T (Y - X\theta) + \delta^2 \theta \overset{\text{identity matrix}}{I} \theta \right\}$$

$$= \frac{\partial}{\partial \theta} \left\{ Y^T Y - 2 Y^T X\theta + \theta^T X^T X\theta + \theta^T (\delta^2 I) \theta \right\}$$

$$= -2 X^T Y + 2 X^T X\theta + 2 \delta^2 I \theta$$
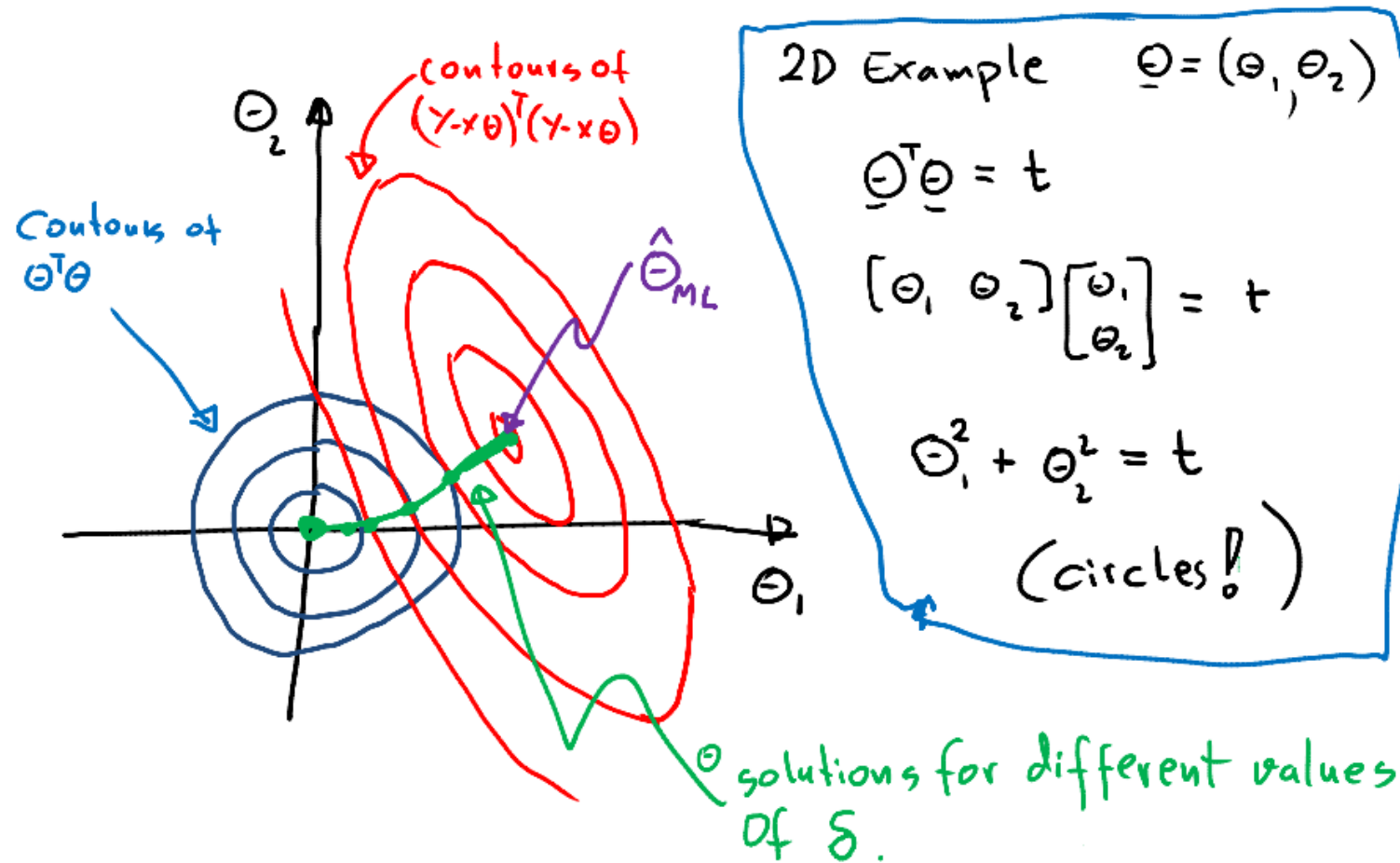
$$= -2 X^T Y + 2 (X^T X + \delta^2 I) \theta$$

Equating to zero, yields
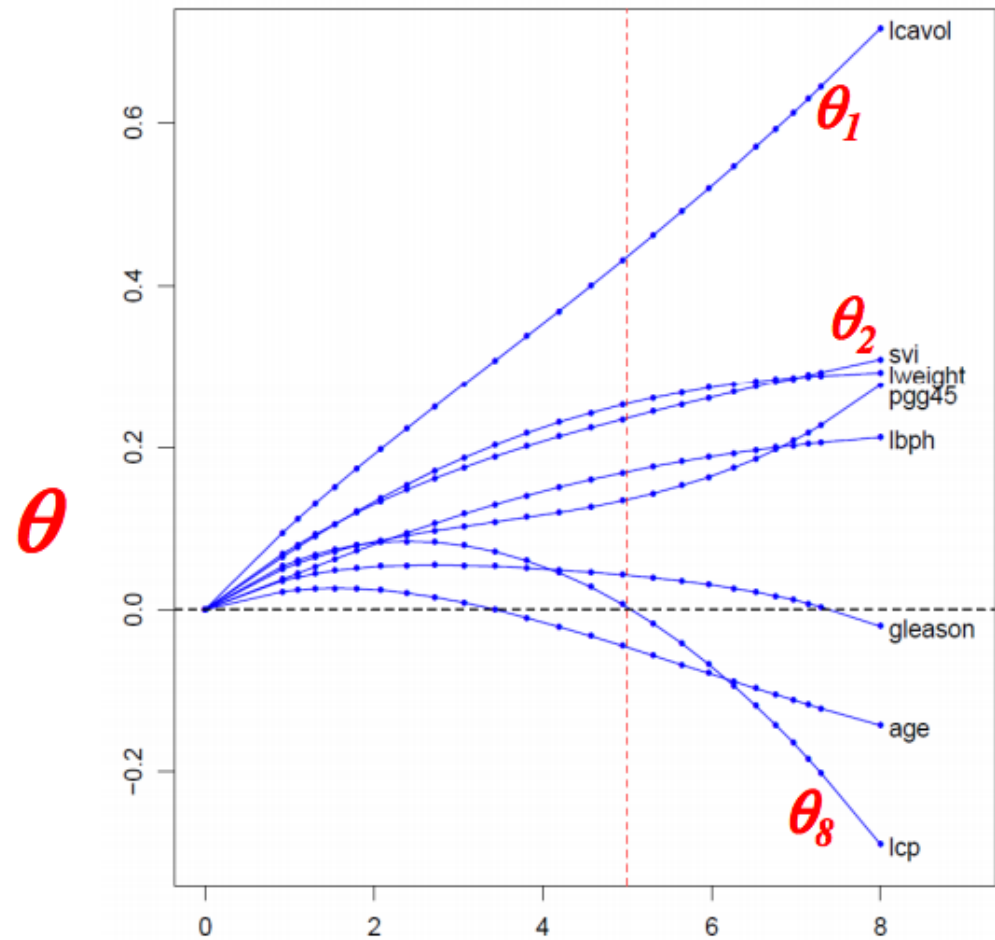
$$\hat{\theta}_{ridge} = (X^T X + \delta^2 I)^{-1} X^T Y$$

# Ridge regression as constrained optimization

$$J(\boldsymbol{\theta}) = (\mathbf{y} - \mathbf{X}\boldsymbol{\theta})^T(\mathbf{y} - \mathbf{X}\boldsymbol{\theta}) + \delta^2\boldsymbol{\theta}^T\boldsymbol{\theta}$$

$$\min_{\boldsymbol{\theta}\,:\,\boldsymbol{\theta}^T\boldsymbol{\theta} \leq t(\delta)} \left\{ (\mathbf{y} - \mathbf{X}\boldsymbol{\theta})^T(\mathbf{y} - \mathbf{X}\boldsymbol{\theta}) \right\}$$



contours of
$(Y - X\Theta)^T (Y - X\Theta)$

Contours of
$\Theta^T \Theta$

$\hat{\Theta}_{ML}$

2D Example    $\underline{\Theta} = (\Theta_1, \Theta_2)$

$\underline{\Theta}^T \underline{\Theta} = t$

$\begin{bmatrix} \Theta_1 & \Theta_2 \end{bmatrix} \begin{bmatrix} \Theta_1 \\ \Theta_2 \end{bmatrix} = t$

$\Theta_1^2 + \Theta_2^2 = t$

(circles!)

$\Theta$ solutions for different values
of $\delta$.

As $\delta$ increases, $t(\delta)$ decreases and each $\theta_i$ goes to zero.

# Ridge, feature selection, shrinkage and weight decay

Large values of $\boldsymbol{\theta}$ are penalised. We are *shrinking* $\boldsymbol{\theta}$ towards zero. This can be used to carry out *feature weighting*. An input $x_{i,d}$ weighted by a small $\theta_d$ will have less influence on the ouptut $y_i$. This penalization with a regularizer is also known as weight decay in the neural networks literature.

Note that shrinking the bias term $\boldsymbol{\theta}_1$ is undesirable. To keep the notation simple, we will assume that the mean of $\mathbf{y}$ has been subtracted from $\mathbf{y}$. This mean is indeed our estimate $\widehat{\boldsymbol{\theta}_1}$.

```python
from keras.regularizers import l2, activity_l2

model.add(Dense(64, input_dim=64, W_regularizer=l2(0.01)))
```

# Selecting features for prediction



$$\hat{y} = x_1 \theta_1^{\,0} + x_2 \theta_2 + \cdots + x_d \theta_d$$

$x_1$ is expensive

$x_1$ does not contribute to good predictions $\hat{y}$

Then we want $\theta_1 \to 0$

# Selecting features for prediction

*As $\delta$ increases, $t(\delta)$ decreases and each $\theta_i$ goes to zero, but too slowly for ridge. Lasso will ensure that irrelevant features $x_i$ have weight $\theta_i = 0$.*



[Hastie, Tibshirani & Friedman book]

# The Lasso: least absolute selection and shrinkage operator

$$J(\theta) = (Y-x\theta)^T (Y-x\theta) + \delta^2 \sum_{j=1}^{d} |\theta_j| \qquad L_1 \text{ Norm}$$

in 2D $\qquad \theta = (\theta_1, \theta_2)$

$$|\theta_1| + |\theta_2| = const$$

$\theta_1 + \theta_2 = const$

$\theta_1 - \theta_2 = const$

$-\theta_1 - \theta_2 = const$

$-\theta_1 + \theta_2 = const$



$\theta_2$

$\hat{\theta}_{ML}$

$\delta^2 = 0$

$\theta_2 = 0$

$\theta_1$

$\delta^2 = \infty$

$J(\theta)$

# Going nonlinear via basis functions

We introduce basis functions $\phi(\cdot)$ to deal with nonlinearity:

$$y(\mathbf{x}) = \phi(\mathbf{x})\boldsymbol{\theta} + \epsilon$$

For example, $\phi(x) = [1, x, x^2]$



$\hat{y} = \phi(x)\Theta$

$= \Theta_0 + x\Theta_1 + x^2\Theta_2$

$\hat{\Theta}_{ML} = \left[\phi(x)^T\phi(x)\right]^{-1}\phi(x)^T y$

# Going nonlinear via basis functions

$$y(\mathbf{x}) = \boldsymbol{\phi}(\mathbf{x})\boldsymbol{\theta} + \epsilon$$

$$\boldsymbol{\phi}(\mathbf{x}) = [1, x_1, x_2]$$

$$x_1 x_2$$

$$\boldsymbol{\phi}(\mathbf{x}) = [1, x_1, x_2, x_1^2, x_2^2]$$

**Example**: *Ridge regression* **with a** *polynomial of degree 14*

$$\hat{y}(x_i) = 1\ \theta_0 + x_i\ \theta_1 + x_i^2\ \theta_2 + \ldots + x_i^{13}\ \theta_{13} + x_i^{14}\ \theta_{14}$$

$$\Phi = [\ 1\ \ x_i\ \ x_i^2\ \ \ldots\ \ x_i^{13}\ \ x_i^{14}\ ]\ \ x_i^{15} \ldots$$

$$J(\theta) = (y - \Phi\theta)^T (y - \Phi\theta) + \delta\ \theta^T \theta$$



small $\delta$

medium $\delta$

large $\delta$

# Kernel regression and RBFs

We can use kernels or radial basis functions (RBFs) as features:

$$\phi(\mathbf{x}) = [\kappa(\mathbf{x}, \boldsymbol{\mu}_1, \lambda), \ldots, \kappa(\mathbf{x}, \boldsymbol{\mu}_d, \lambda)], \quad e.g. \quad \kappa(\mathbf{x}, \boldsymbol{\mu}_i, \lambda) = e^{(-\frac{1}{\lambda}\|\mathbf{x} - \boldsymbol{\mu}_i\|^2)}$$

We can choose the locations $\mu$ of the **basis functions** to be the inputs.
That is, $\mu_i = x_i$. These basis functions are the known as **kernels**.
The choice of width $\lambda$ is tricky, as illustrated below.

**kernels**



*Too small $\lambda$*

*Right $\lambda$*

*Too large $\lambda$*

The big question is how do we choose the regularization coefficient, the width of the kernels or the polynomial order?
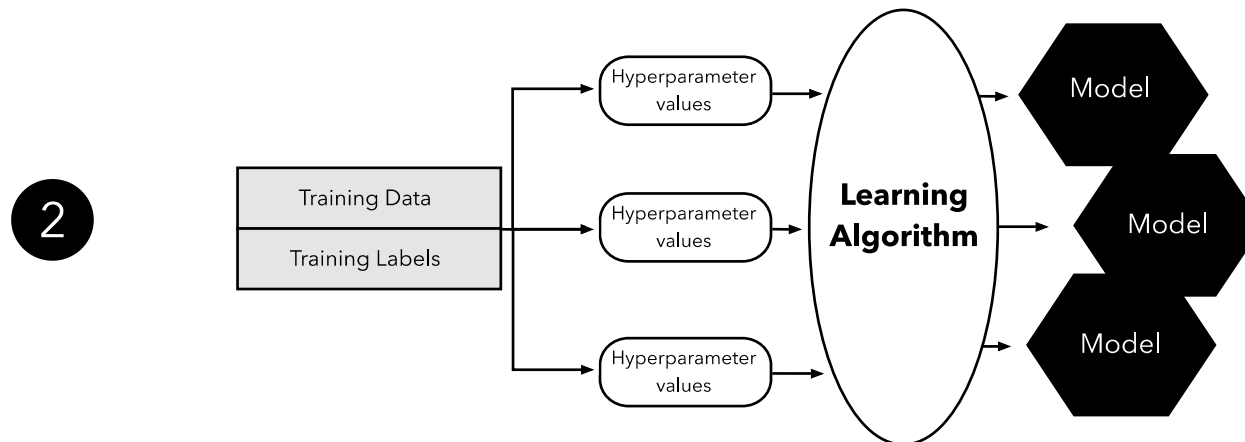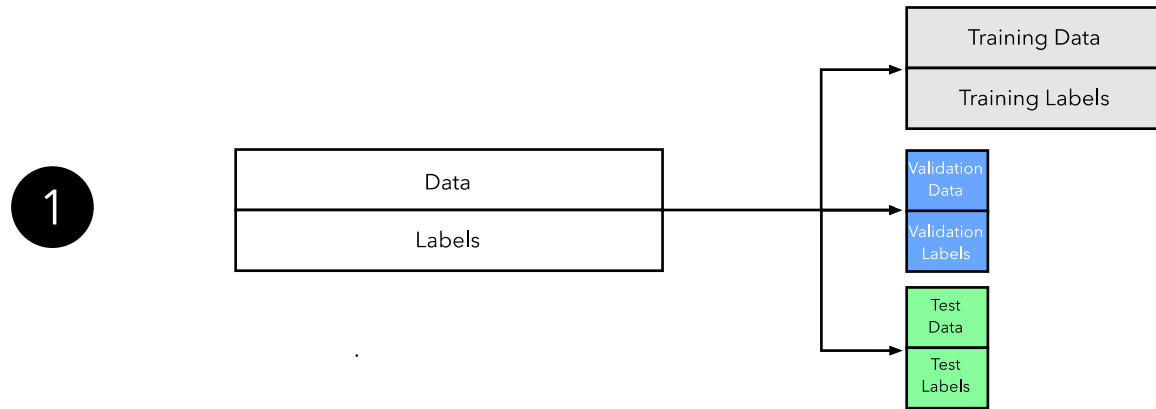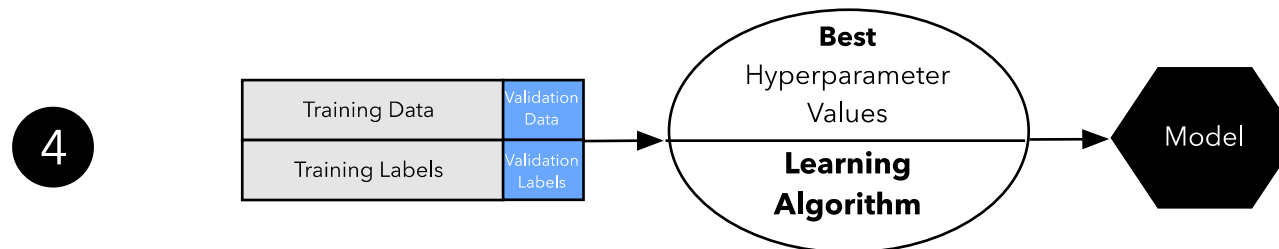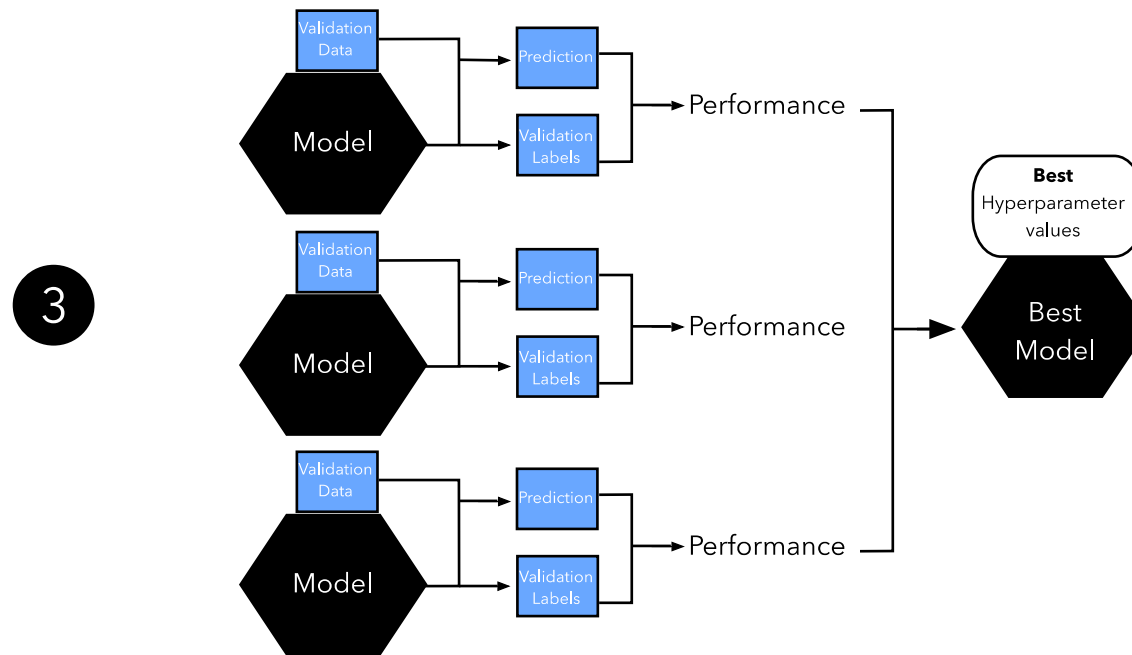
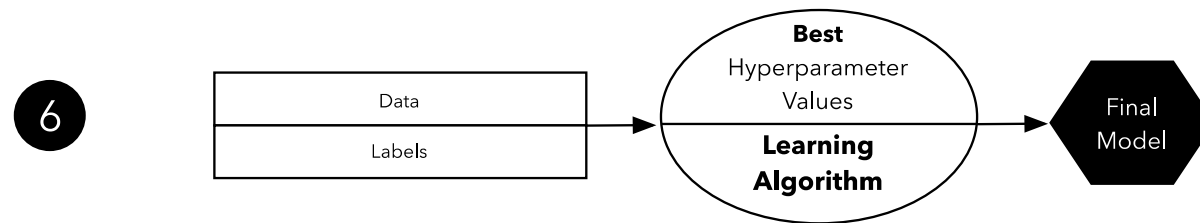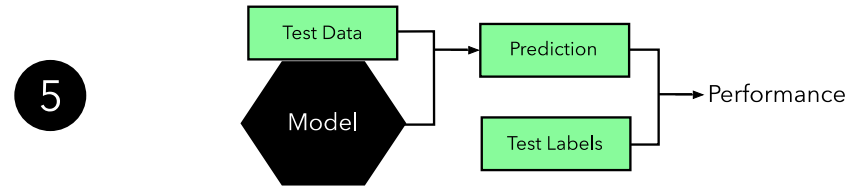# Holdout Evaluation I

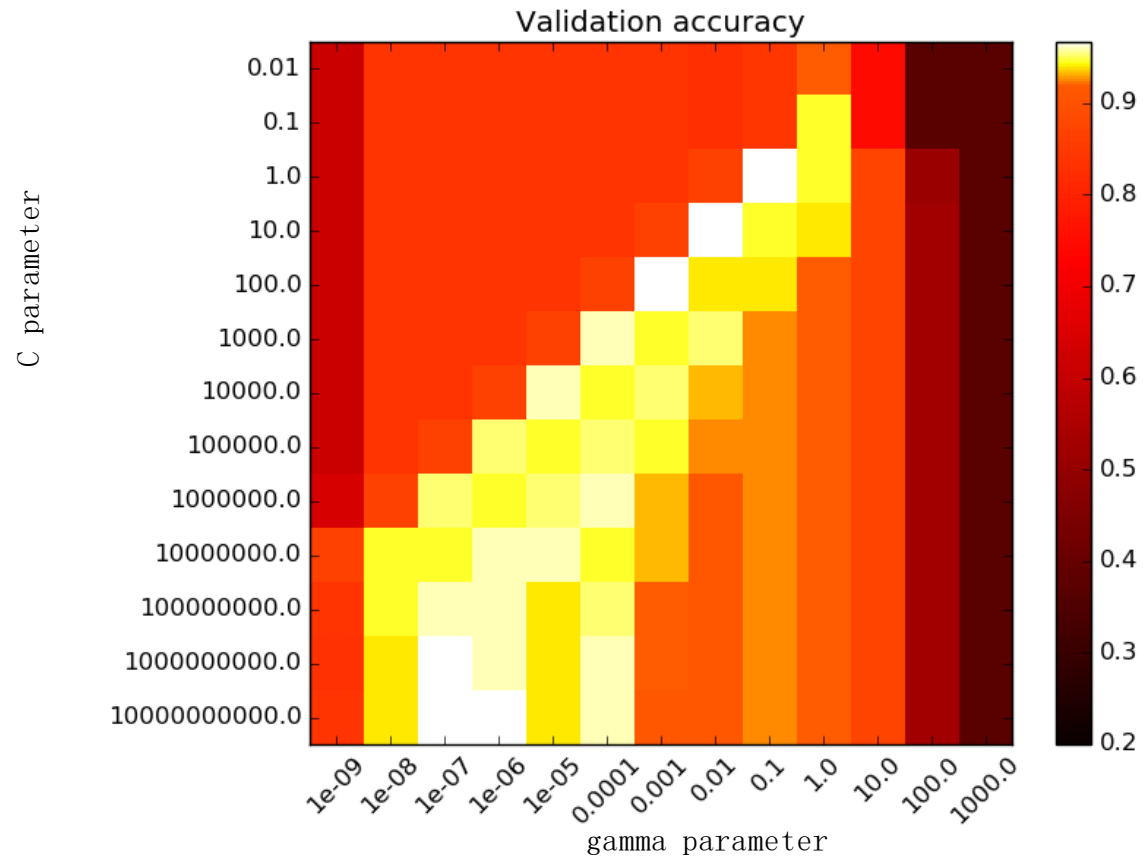# Holdout Evaluation II

**3**



**4**

# Holdout Validation I

# Holdout Validation II

# Holdout Validation III

# Grid Search

# Now, big question

- How to define input X?

- http://stockcharts.com/school/doku.php?id=chart_school:technical_indicators