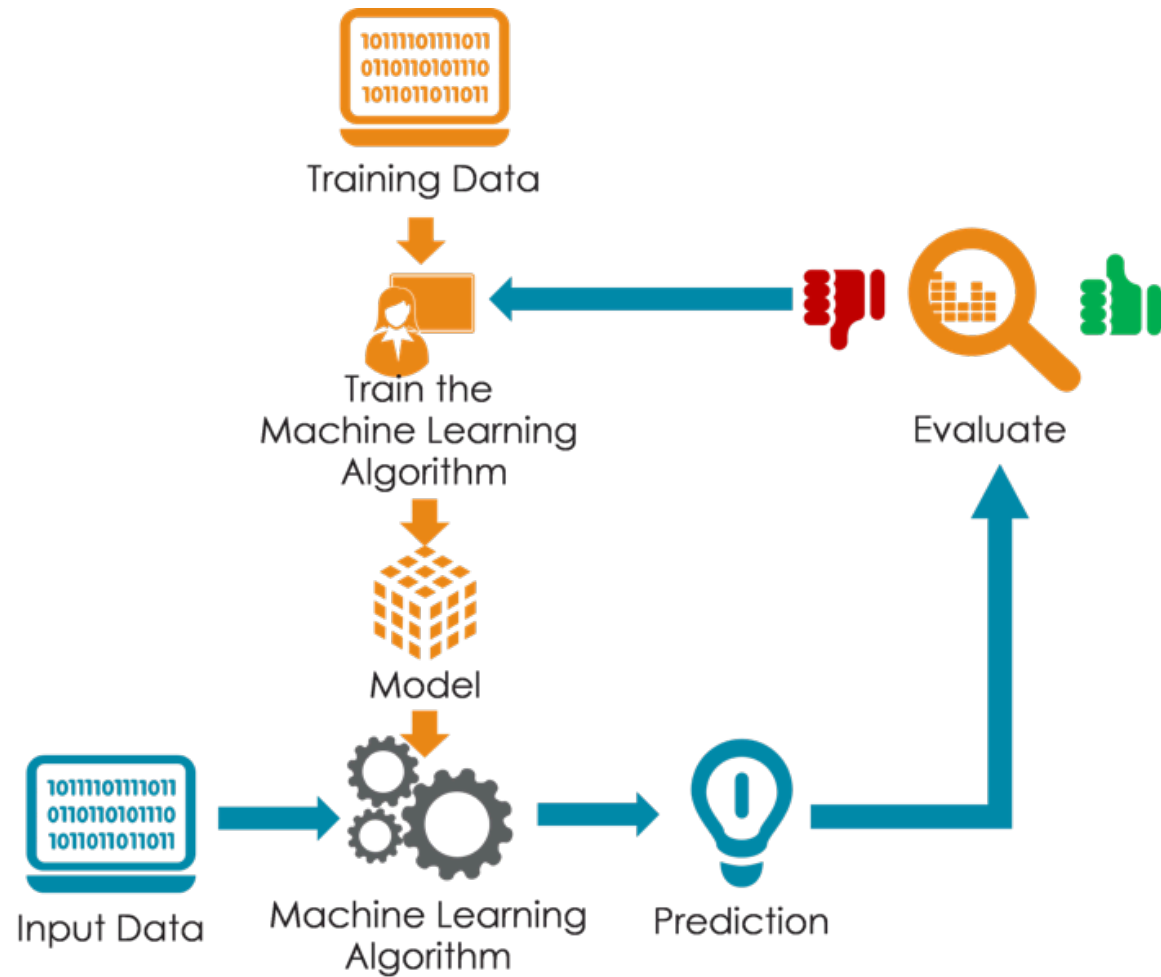


MACHINE LEARNING



SUPERVISED LEARNING

Predict **outputs** based on **inputs**

SUPERVISED LEARNING

independent variables
predictors
attributes
features

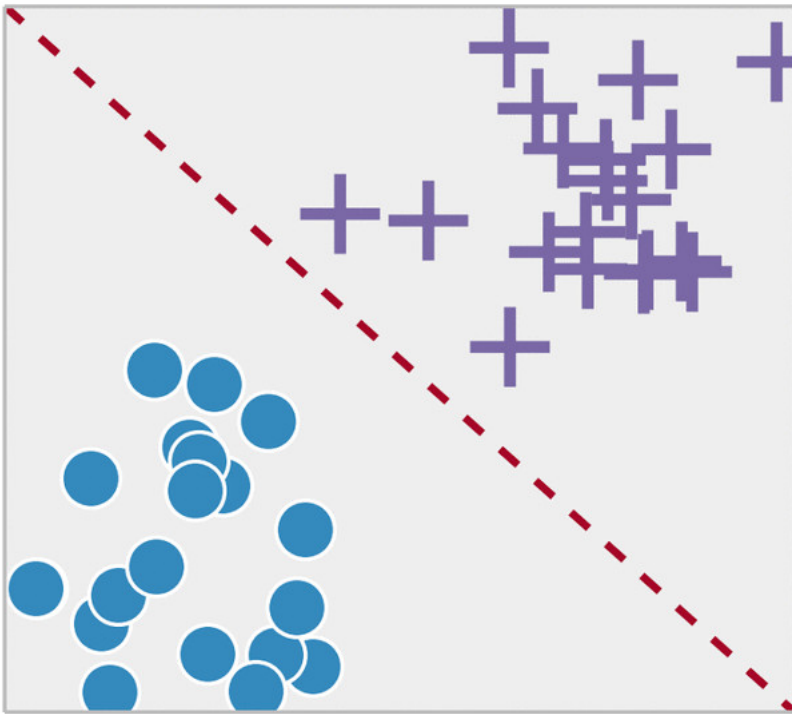
Predict **outputs** based on **inputs**

targets
responses
outcomes
dependent variables

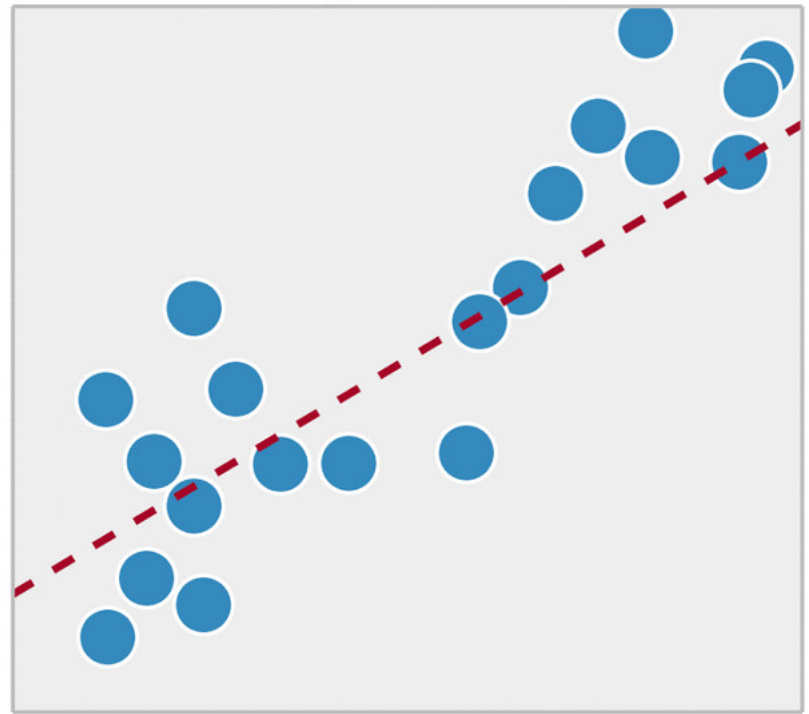
SUPERVISED LEARNING

Predict **outputs** based on **inputs**

Classification



Regression



UNSUPERVISED LEARNING

There is no **output** measure

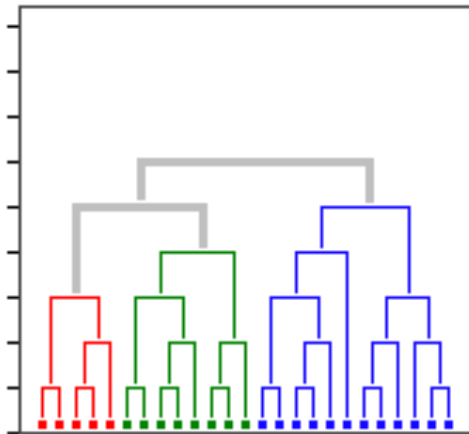
Describe associations/patterns among **inputs**

UNSUPERVISED LEARNING

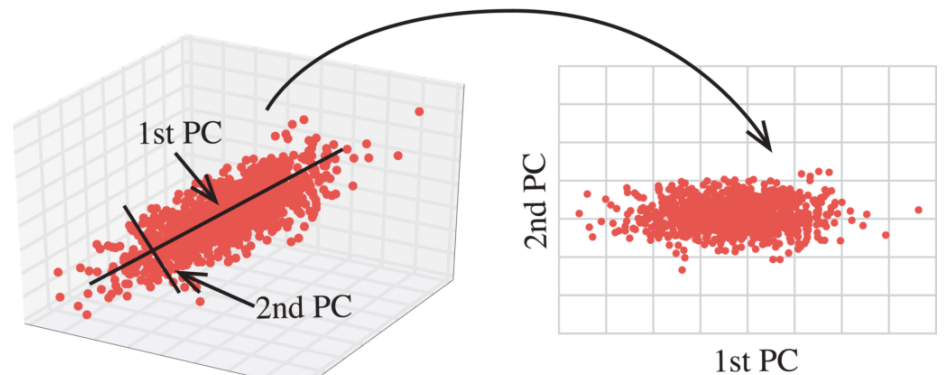
There is no **output** measure

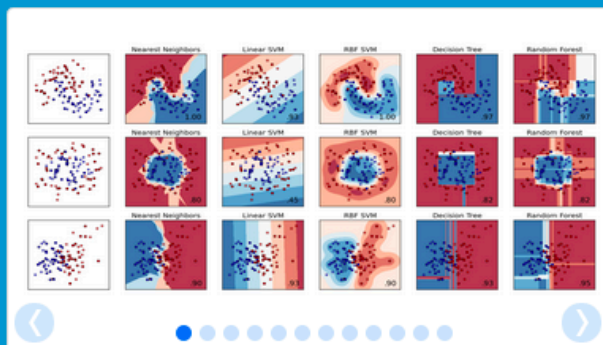
Describe associations/patterns among **inputs**

Clustering



Dimensionality Reduction





scikit-learn

Machine Learning in Python

- Simple and efficient tools for data mining and data analysis
- Accessible to everybody, and reusable in various contexts
- Built on NumPy, SciPy, and matplotlib
- Open source, commercially usable - BSD license

Classification

Identifying to which set of categories a new observation belong to.

Applications: Spam detection, Image recognition.

Algorithms: *SVM, nearest neighbors, random forest, ...* — [Examples](#)

Regression

Predicting a continuous value for a new example.

Applications: Drug response, Stock prices.

Algorithms: *SVR, ridge regression, Lasso, ...* — [Examples](#)

Clustering

Automatic grouping of similar objects into sets.

Applications: Customer segmentation, Grouping experiment outcomes

Algorithms: *k-Means, spectral clustering, mean-shift, ...* — [Examples](#)

Dimensionality reduction

Reducing the number of random variables to consider.

Applications: Visualization, Increased efficiency

Algorithms: *PCA, Isomap, non-negative matrix factorization.* — [Examples](#)

Model selection

Comparing, validating and choosing parameters and models.

Goal: Improved accuracy via parameter tuning

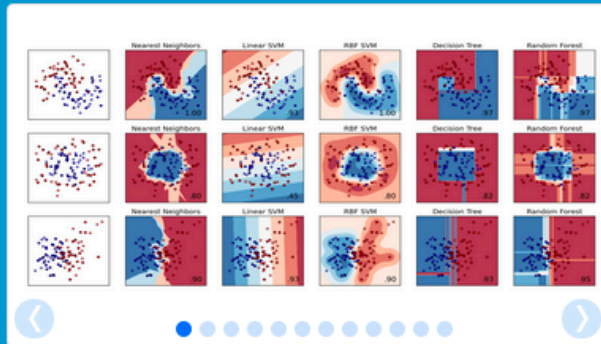
Modules: *grid search, cross validation, metrics.* — [Examples](#)

Preprocessing

Feature extraction and normalization.

Application: Transforming input data such as text for use with machine learning algorithms.

Modules: *preprocessing, feature extraction.* — [Examples](#)



scikit-learn

Machine Learning in Python

- Simple and efficient tools for data mining and data analysis
- Accessible to everybody, and reusable in various contexts
- Built on NumPy, SciPy, and matplotlib
- Open source, commercially usable - BSD license

Classification

Identifying to which set of categories a new observation belong to.

Applications: Spam detection, Image recognition.

Algorithms: *SVM, nearest neighbors, random forest, ...* — [Examples](#)

Regression

Predicting a continuous value for a new example.

Applications: Drug response, Stock prices.

Algorithms: *SVR, ridge regression, Lasso, ...* — [Examples](#)

Clustering

Automatic grouping of similar objects into sets.

Applications: Customer segmentation, Grouping experiment outcomes

Algorithms: *k-Means, spectral clustering, mean-shift, ...* — [Examples](#)

Dimensionality reduction

Reducing the number of random variables to consider.

Applications: Visualization, Increased efficiency

Algorithms: *PCA, Isomap, non-negative matrix factorization.* — [Examples](#)

Model selection

Comparing, validating and choosing parameters and models.

Goal: Improved accuracy via parameter tuning

Modules: *grid search, cross validation, metrics.* — [Examples](#)

Preprocessing

Feature extraction and normalization.

Application: Transforming input data such as text for use with machine learning algorithms.

Modules: *preprocessing, feature extraction.* — [Examples](#)

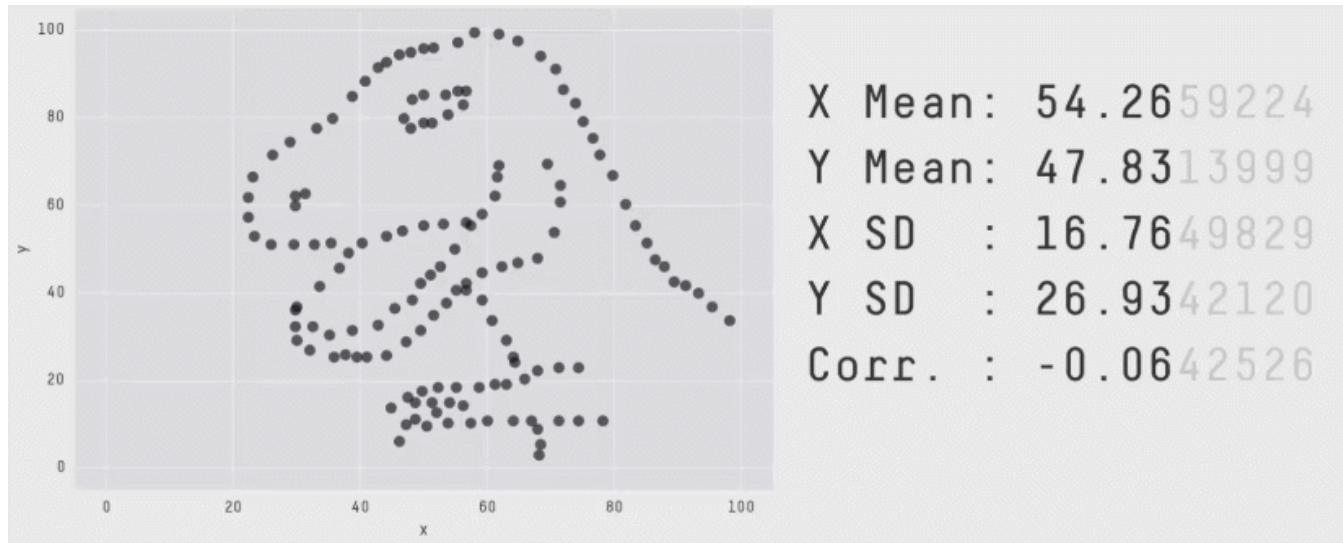
CORRELATION

- Measure the direction and strength of a linear relationship.

$$r = \frac{\sum_{i=1}^n ((x_i - \bar{x})(y_i - \bar{y}))}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2 \sum_{i=1}^n (y_i - \bar{y})^2}}$$

CORRELATION

- Measure the direction and strength of a **linear relationship**.



CORRELATION

- Measure the direction and strength of a linear relationship.
- Cannot be used for prediction purposes.

REGRESSION ANALYSIS

1. State the problem
2. Select potentially relevant variables
3. Specify the model
4. Fit the model
5. Critically evaluate the model
6. Address the original problem

LINEAR REGRESSION

- Express functional relationships among variables as an equation or “model”

$$Y = \beta_0 + \beta_1 X + \varepsilon$$

LINEAR REGRESSION

- Express functional relationships among variables as an equation or “model”

The diagram illustrates the linear regression equation $Y = \beta_0 + \beta_1 X + \epsilon$ with the following components and labels:

- Y**: response variable (indicated by an orange box and a downward arrow)
- β_0** : intercept (indicated by a blue box and an upward arrow labeled "parameters to be estimated 'coefficients'" and "intercept")
- β_1** : slope (indicated by a blue box and an upward arrow labeled "parameters to be estimated 'coefficients'" and "slope")
- X**: predictor variable (indicated by a yellow box and a downward arrow)
- ϵ** : random error "noise" (indicated by a green box and an upward arrow labeled "random error 'noise'" and "noise")

The equation is presented as $Y = \beta_0 + \beta_1 X + \epsilon$.

LINEAR REGRESSION

Simple linear regression:

$$Y = \beta_0 + \beta_1 X + \varepsilon$$

Multiple linear regression:

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_p X_p + \varepsilon$$

LINEAR REGRESSION

Linear refers to the way that the parameters enter the model...

Linear

$$Y = \beta_0 + \beta_1 X_1 + \varepsilon$$

$$Y = \beta_0 + \beta_1 X + \beta_2 X^2 + \varepsilon$$

Non-linear

$$Y = \beta_0 + e^{\beta_1 X_1} + \varepsilon$$

DATA NOTATION

Observation Number	Response Y	Predictors			
		X_1	X_2	\dots	X_p
1	y_1	x_{11}	x_{12}	\dots	x_{1p}
2	y_2	x_{21}	x_{22}	\dots	x_{2p}
3	y_3	x_{31}	x_{32}	\dots	x_{3p}
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots
n	y_n	x_{n1}	x_{n2}	\dots	x_{np}

DATA NOTATION

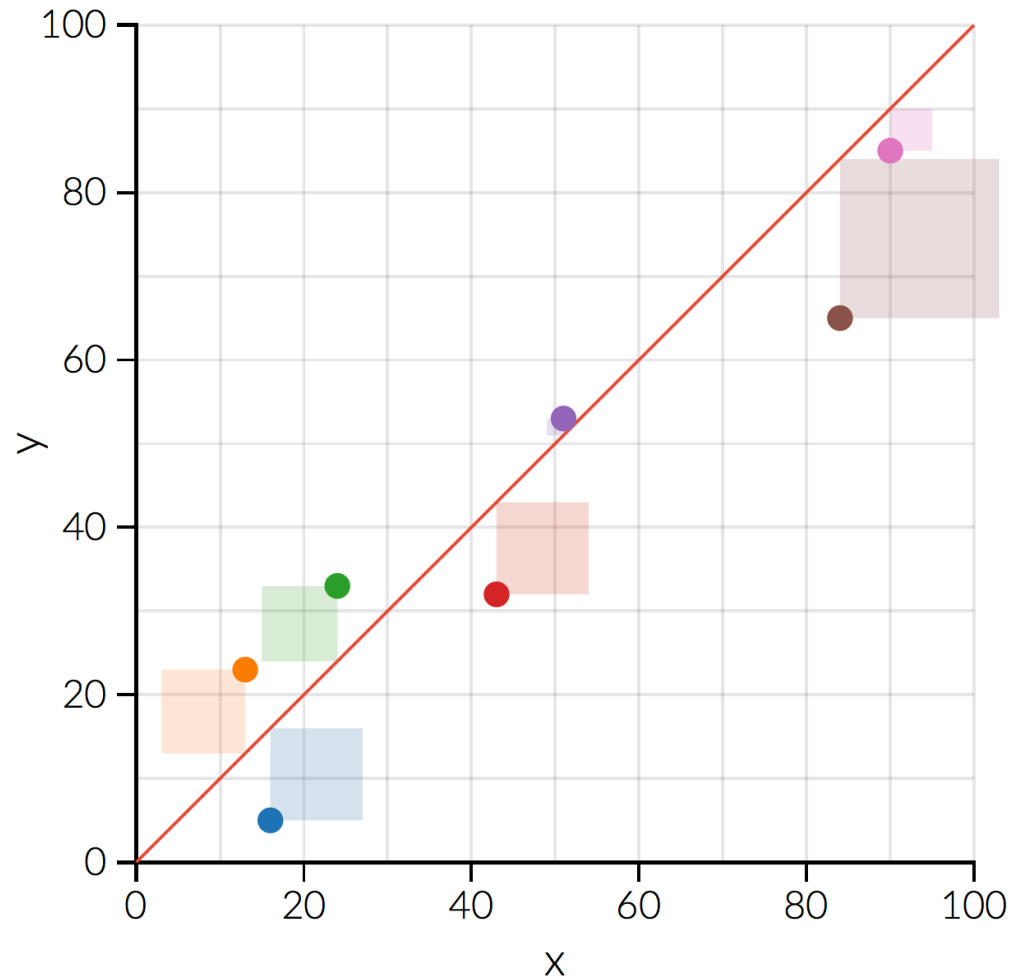
Observation Number	Response Y	Predictors			
		X_1	X_2	\dots	X_p
1	y_1	x_{11}	x_{12}	\dots	x_{1p}
2	y_2	x_{21}	x_{22}	\dots	x_{2p}
3	y_3	x_{31}	x_{32}	\dots	x_{3p}
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots
n	y_n	x_{n1}	x_{n2}	\dots	x_{np}

...what does this data structure resemble?

MODEL FITTING

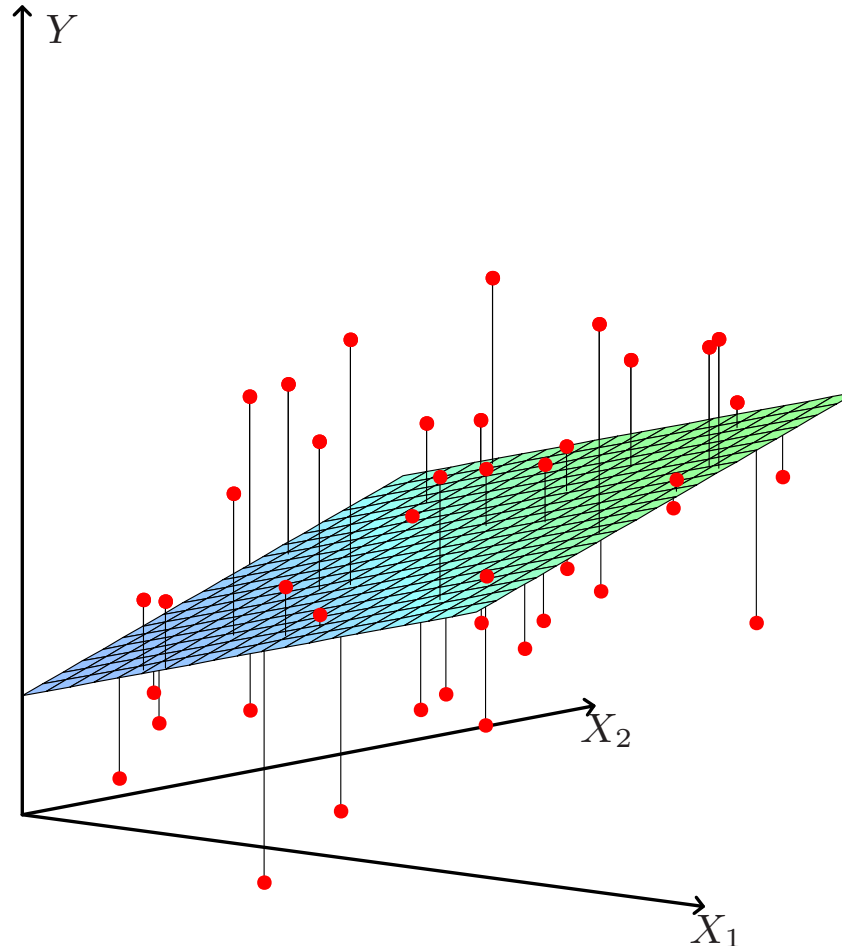
Minimize the sum of squared residuals...

$$\sum_{i=1}^n (y_i - \hat{y}_i)^2$$



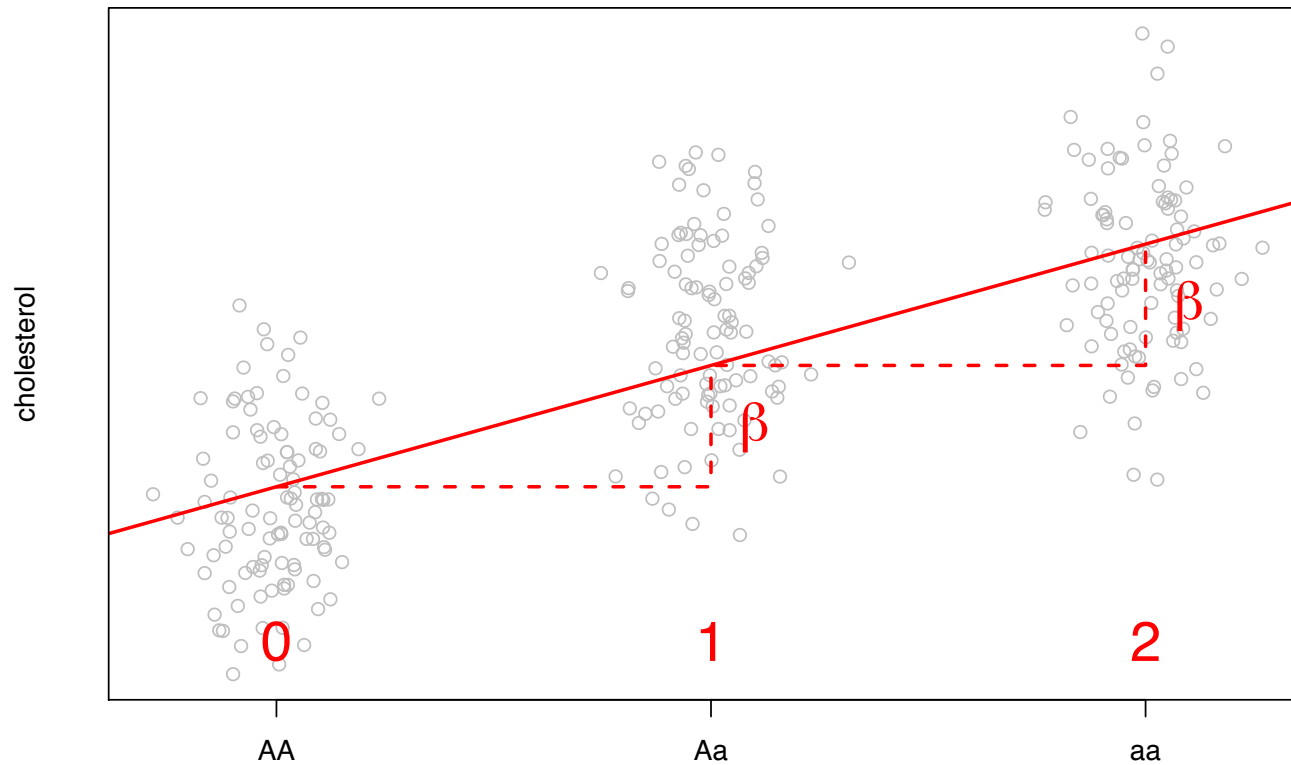
MODEL FITTING

Same concept applies in two or more dimensions...



APPLICATIONS OF LINEAR REGRESSION

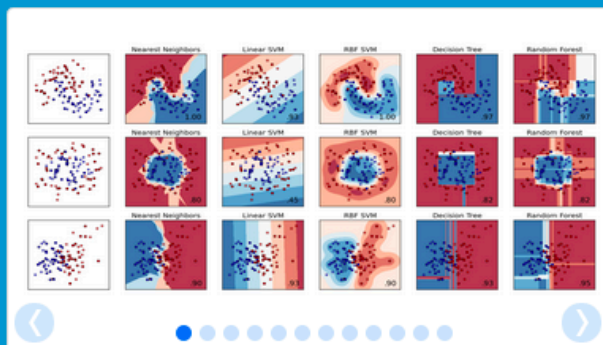
$$y = \beta_0 + \beta \times \# \text{minor alleles}$$



GENOME-WIDE ASSOCIATION STUDIES

Hundreds of thousands of linear regressions...





scikit-learn

Machine Learning in Python

- Simple and efficient tools for data mining and data analysis
- Accessible to everybody, and reusable in various contexts
- Built on NumPy, SciPy, and matplotlib
- Open source, commercially usable - BSD license

Classification

Identifying to which set of categories a new observation belong to.

Applications: Spam detection, Image recognition.

Algorithms: *SVM, nearest neighbors, random forest, ...* — Examples

Regression

Predicting a continuous value for a new example.

Applications: Drug response, Stock prices.

Algorithms: *SVR, ridge regression, Lasso, ...* — Examples

Clustering

Automatic grouping of similar objects into sets.

Applications: Customer segmentation, Grouping experiment outcomes

Algorithms: *k-Means, spectral clustering, mean-shift, ...* — Examples

Dimensionality reduction

Reducing the number of random variables to consider.

Applications: Visualization, Increased efficiency

Algorithms: *PCA, Isomap, non-negative matrix factorization.* — Examples

Model selection

Comparing, validating and choosing parameters and models.

Goal: Improved accuracy via parameter tuning

Modules: *grid search, cross validation, metrics.* — Examples

Preprocessing

Feature extraction and normalization.

Application: Transforming input data such as text for use with machine learning algorithms.

Modules: *preprocessing, feature extraction.* — Examples



scikit-learn

Machine Learning in Python

- Simple and efficient tools for data mining and data analysis
- Accessible to everybody, and reusable in various contexts
- Built on NumPy, SciPy, and matplotlib
- Open source, commercially usable - BSD license

Classification

Identifying to which set of categories a new observation belong to.

Applications: Spam detection, Image recognition.

Algorithms: *SVM, nearest neighbors, random forest, ...* — Examples

Regression

Predicting a continuous value for a new example.

Applications: Drug response, Stock prices.

Algorithms: *SVR, ridge regression, Lasso, ...* — Examples

Clustering

Automatic grouping of similar objects into sets.

Applications: Customer segmentation, Grouping experiment outcomes

Algorithms: *k-Means, spectral clustering, mean-shift, ...* — Examples

Dimensionality reduction

Reducing the number of random variables to consider.

Applications: Visualization, Increased efficiency

Algorithms: *PCA, Isomap, non-negative matrix factorization.* — Examples

Model selection

Comparing, validating and choosing parameters and models.

Goal: Improved accuracy via parameter tuning

Modules: *grid search, cross validation, metrics.* — Examples

Preprocessing

Feature extraction and normalization.

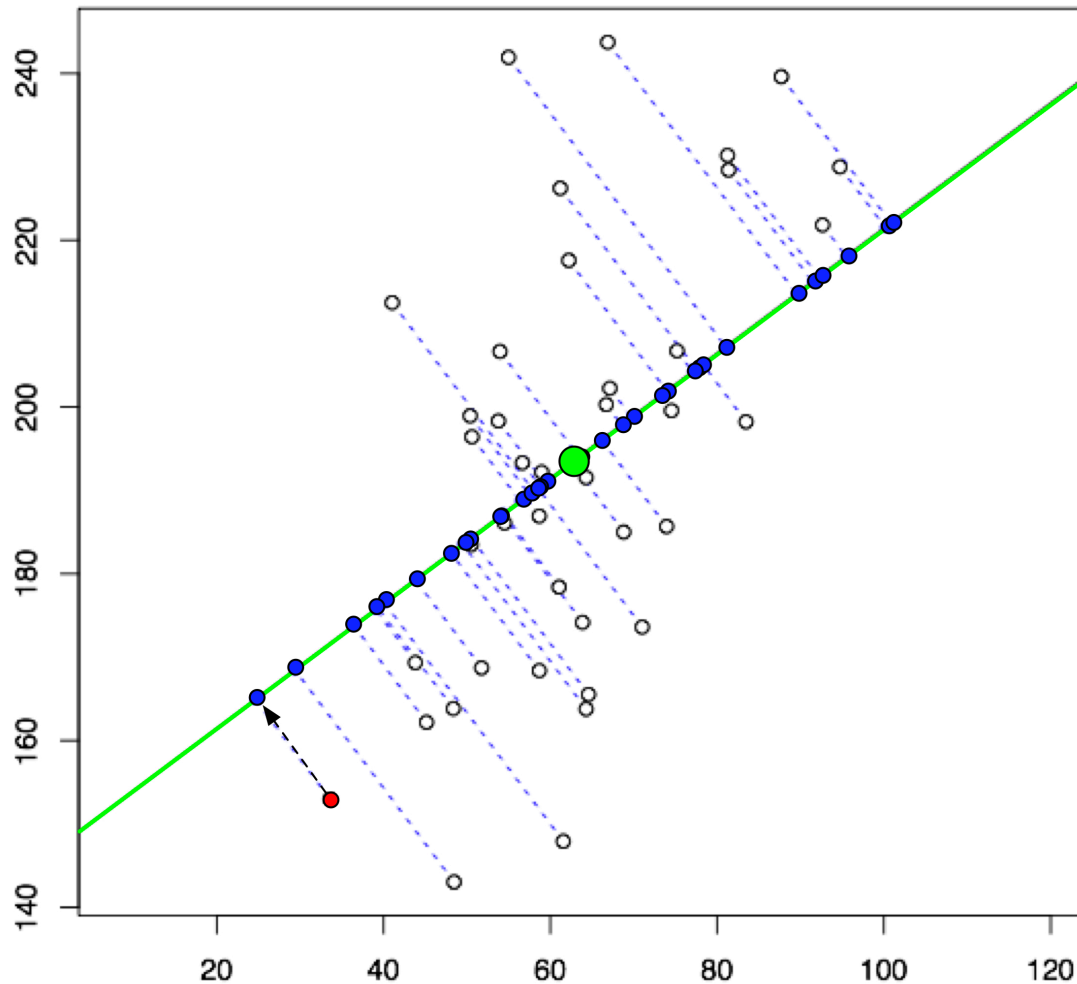
Application: Transforming input data such as text for use with machine learning algorithms.

Modules: *preprocessing, feature extraction.* — Examples

PRINCIPAL COMPONENT ANALYSIS

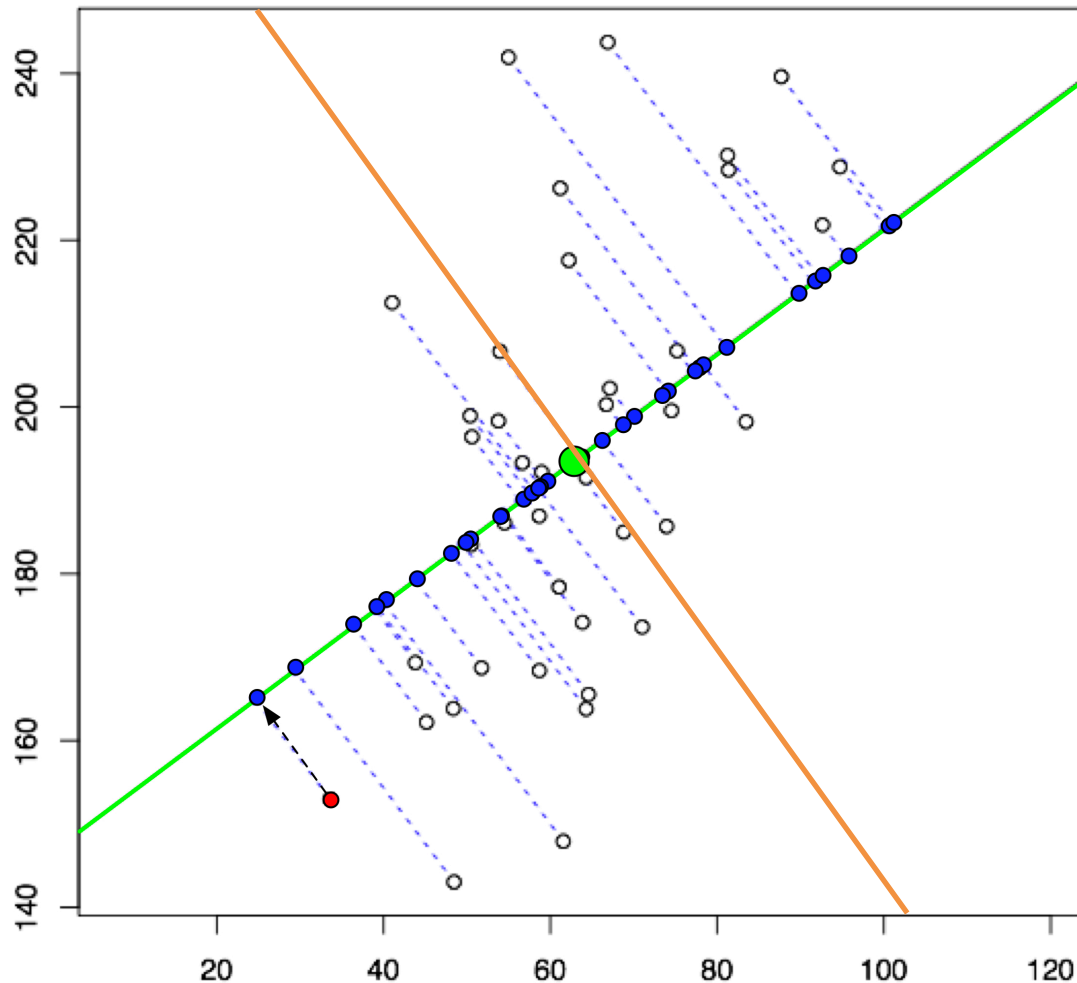
Transform a large number of possibly correlated variables into a smaller number of uncorrelated variables (“principal components”) that capture as much information as possible from the original dataset.

PRINCIPAL COMPONENT ANALYSIS



The **first principal component** is a linear combination of original predictor variables which captures the maximum variance in the data set.

PRINCIPAL COMPONENT ANALYSIS

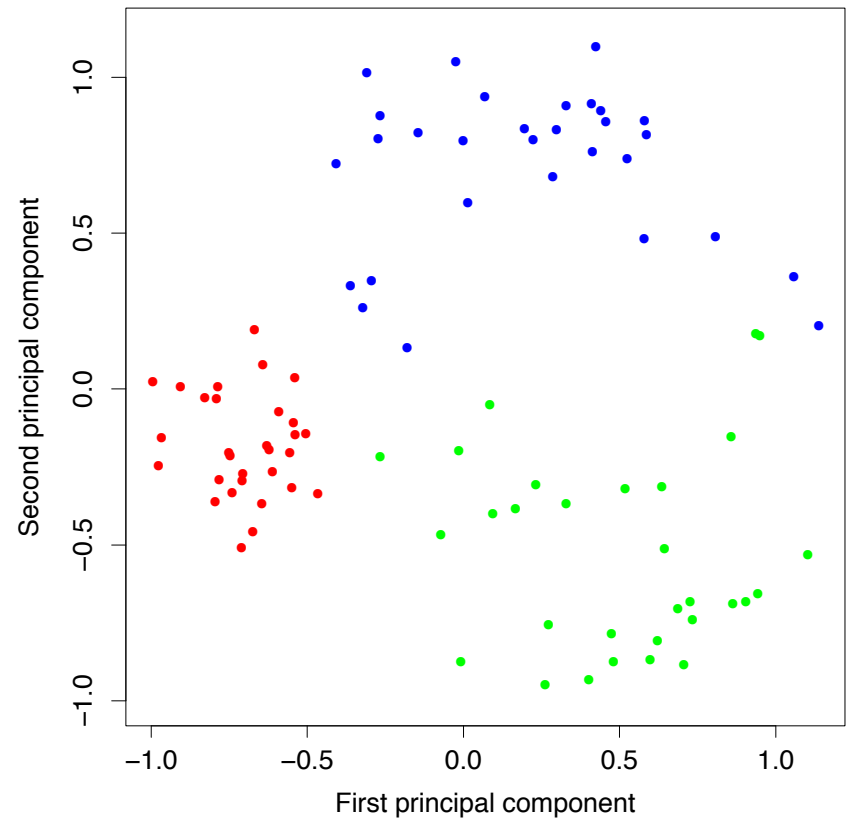
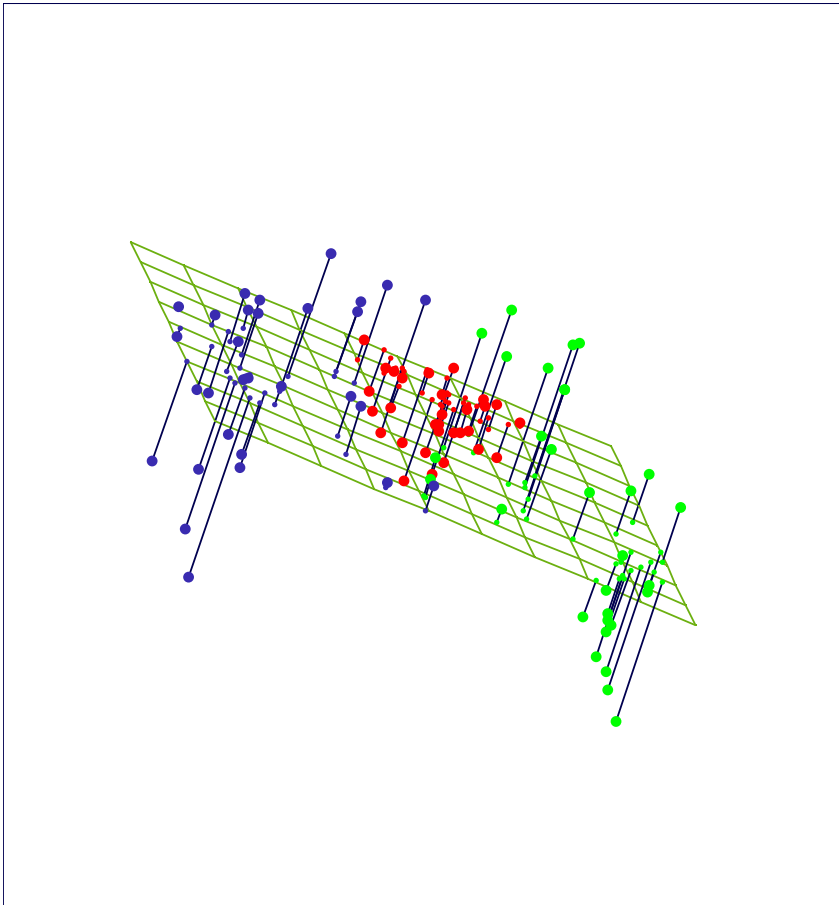


The **first principal component** is a linear combination of original predictor variables which captures the maximum variance in the data set.

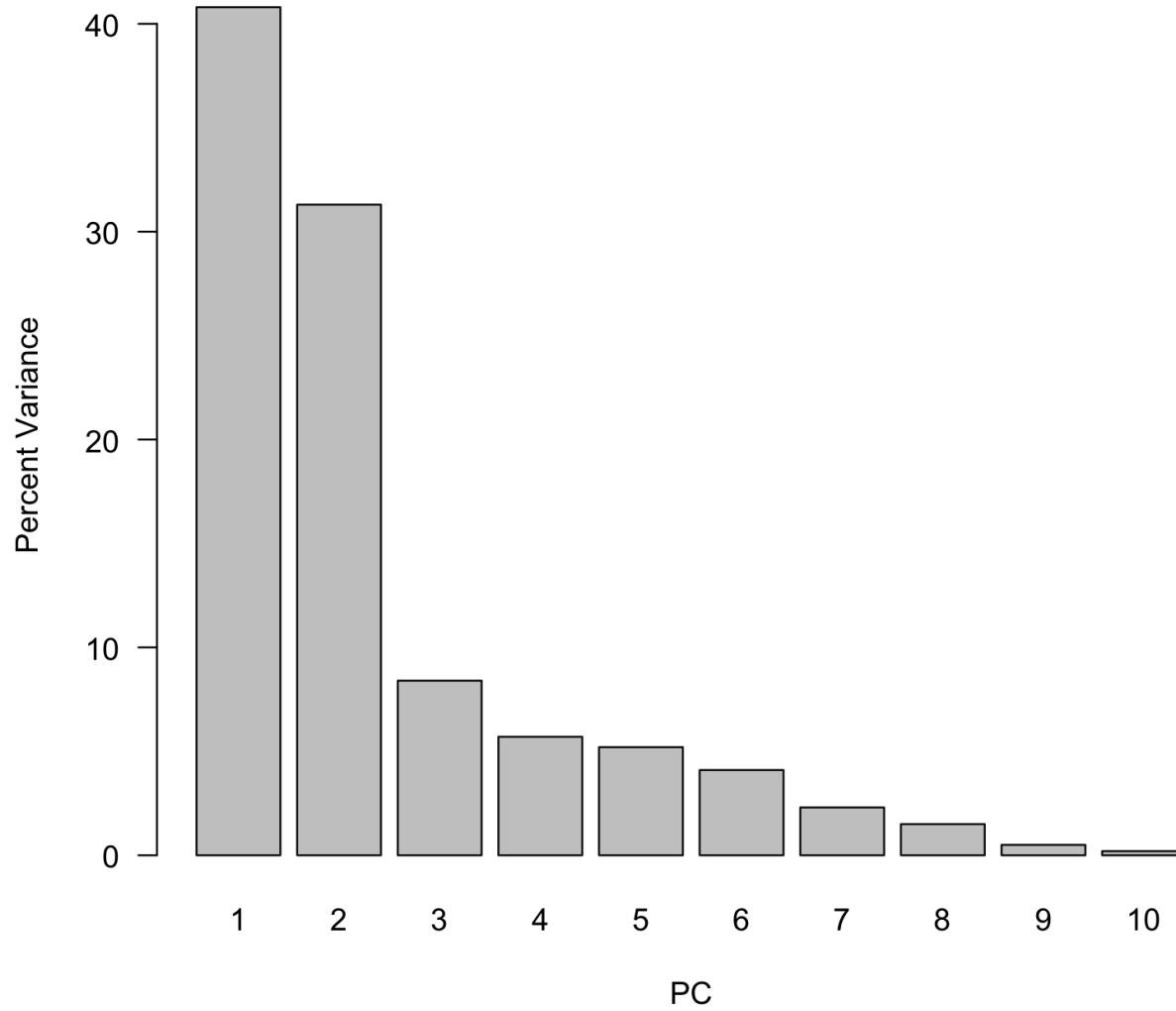
The **second principal component** is a linear combination of original predictor variables which captures the second most variance under the constraint that it is orthogonal to **PC1**.

PRINCIPAL COMPONENT ANALYSIS

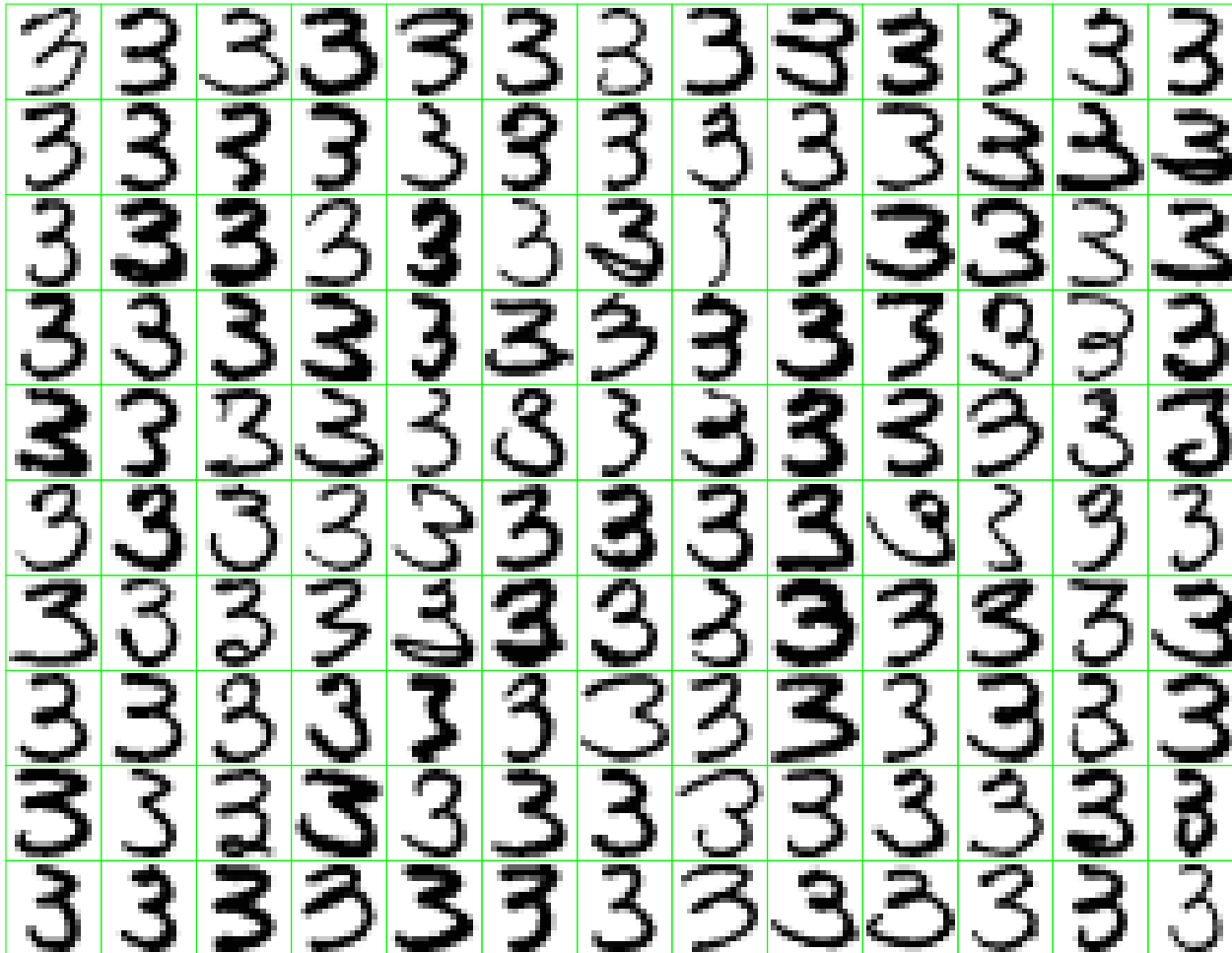
Reduce dimensionality and visualize data structure...



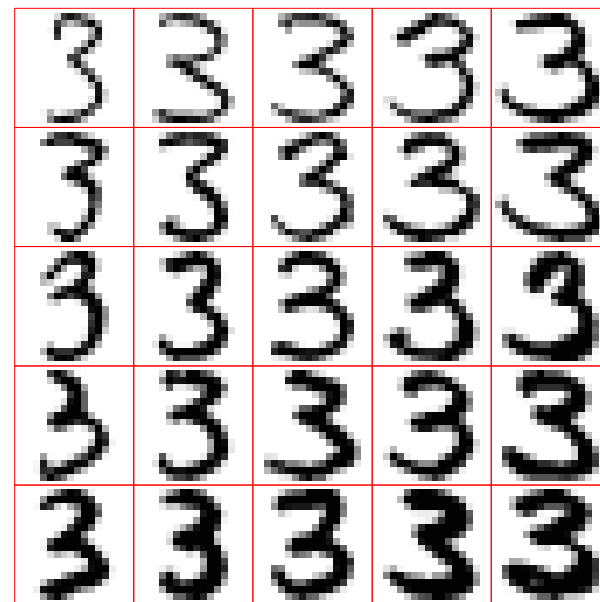
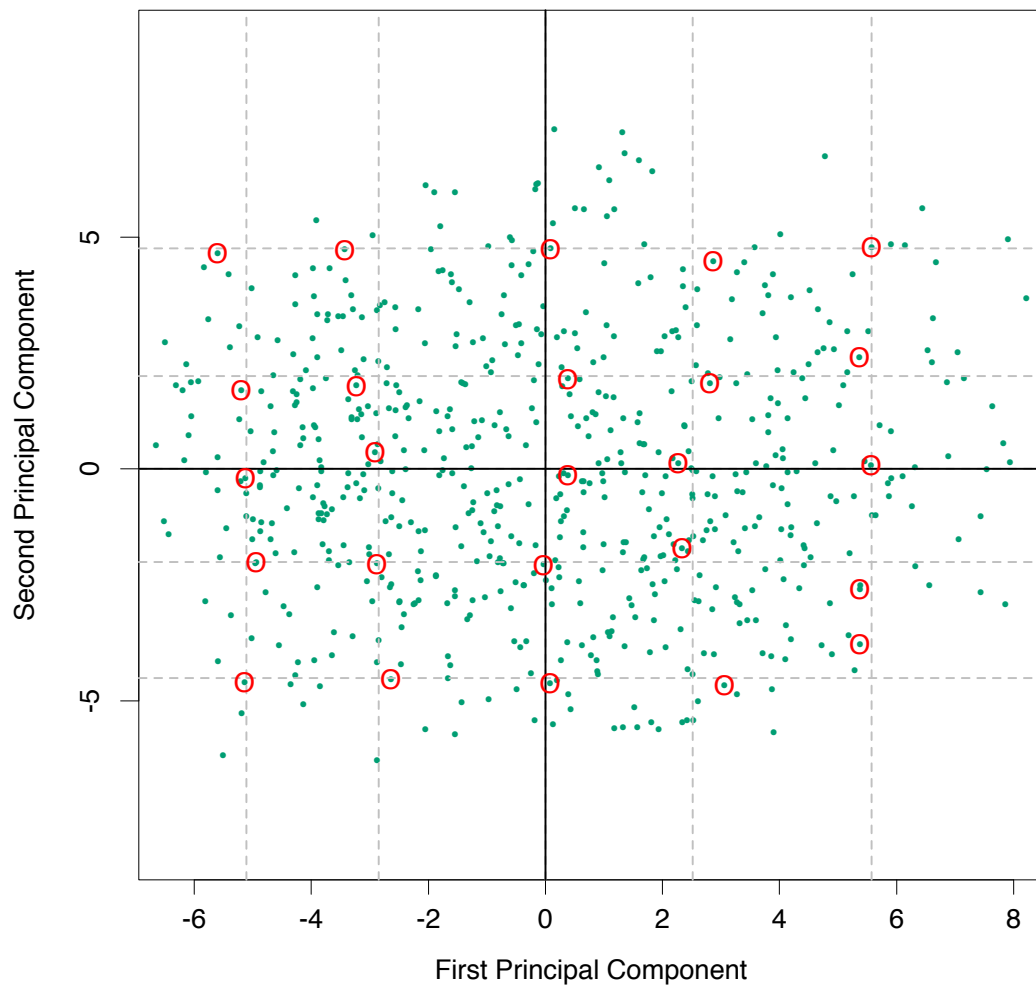
PRINCIPAL COMPONENT ANALYSIS



APPLICATIONS OF PCA



APPLICATIONS OF PCA



APPLICATIONS OF PCA

	rs10001	rs10002	rs10003	rs10004	rs10005	rs10006	...
Sample1	AA	GG	CC	AC	GG	CC	
Sample2	AA	GG	CT	AC	CG	CC	
Sample3	AT	GG	TT	AA	CG	TT	
Sample4	TT	GG	TT	AA	GG	CC	
Sample5	AA	GG	TT	AA	GG	TT	
Sample6	AA	GA	CC	AC	GG	CC	
Sample7	AT	GG	CC	CC	GG	CC	
Sample8	AA	GG	CT	AA	GG	CT	

...

genotypes can
also be represented
counts of non-reference
alleles (i.e., 0, 1, 2)

$$\begin{bmatrix} 0 & 2 & 1 & 2 & 0 & 2 & 2 & 2 \\ 2 & 0 & 0 & 2 & 1 & 0 & 2 & 2 \\ 0 & 1 & 0 & 2 & 0 & 0 & 2 & 0 \\ 1 & 2 & 2 & 0 & 2 & 1 & 0 & 0 \end{bmatrix}$$

APPLICATIONS OF PCA

