Would you like to see this cheatsheet in your native language? You can help us **translating it** (https://github.com/shervine/cheatsheet-translation) on GitHub!

**(https://stanford.edu/~shervine/teaching/cs-221/cheatsheet-reflex-models#cs-221---artificial-intelligence)CS 221 - Artificial Intelligence (teaching/cs-221)** | English ▾ | ⊕

| Reflex | States | Variables | Logic |

# (https://stanford.edu/~shervine/teaching/cs-221/cheatsheet-reflex-models#cheatsheet)Reflex-based models with Machine Learning

*By Afshine Amidi (https://twitter.com/afshinea) and Shervine Amidi (https://twitter.com/shervine)*

Star (https://github.com/afshinea/stanford-cs-221-artificial-intelligence)

**(https://stanford.edu/~shervine/teaching/cs-221/cheatsheet-reflex-models#linear-predictors)**
## Linear predictors

In this section, we will go through reflex-based models that can improve with experience, by going through samples that have input-output pairs.

❏ **Feature vector** — The feature vector of an input $x$ is denoted $\phi(x)$ and is such that:

$$\phi(x) = \begin{bmatrix} \phi_1(x) \\ \vdots \\ \phi_d(x) \end{bmatrix} \in \mathbb{R}^d$$

❏ **Score** — The score $s(x, w)$ of an example $(\phi(x), y) \in \mathbb{R}^d \times \mathbb{R}$ associated to a linear model of weights $w \in \mathbb{R}^d$ is given by the inner product:

$$s(x, w) = w \cdot \phi(x)$$

### Classification

❏ **Linear classifier** — Given a weight vector $w \in \mathbb{R}^d$ and a feature vector $\phi(x) \in \mathbb{R}^d$, the binary linear classifier $f_w$ is given by:

$$f_w(x) = \text{sign}(s(x, w)) = \begin{cases} +1 & \text{if } w \cdot \phi(x) > 0 \\ -1 & \text{if } w \cdot \phi(x) < 0 \\ ? & \text{if } w \cdot \phi(x) = 0 \end{cases}$$



❏ **Margin** — The margin $m(x, y, w) \in \mathbb{R}$ of an example $(\phi(x), y) \in \mathbb{R}^d \times \{-1, +1\}$ associated to a linear model of weights $w \in \mathbb{R}^d$ quantifies the confidence of the prediction: larger values are better. It is given by:

$$m(x, y, w) = s(x, w) \times y$$

### Regression

❏ **Linear regression** — Given a weight vector $w \in \mathbb{R}^d$ and a feature vector $\phi(x) \in \mathbb{R}^d$, the output of a linear regression of weights $w$ denoted as $f_w$ is given by:

$$f_w(x) = s(x, w)$$

❏ **Residual** — The residual $\text{res}(x, y, w) \in \mathbb{R}$ is defined as being the amount by which the prediction $f_w(x)$ overshoots the target $y$:

$$\text{res}(x, y, w) = f_w(x) - y$$

**(https://stanford.edu/~shervine/teaching/cs-221/cheatsheet-reflex-models#loss-minimization)**
## Loss minimization

❏ **Loss function** — A loss function $\text{Loss}(x, y, w)$ quantifies how unhappy we are with the weights $w$ of the model in the prediction task of output $y$ from input $x$. It is a quantity we want to minimize during the training process.

❏ **Classification case** — The classification of a sample $x$ of true label $y \in \{-1, +1\}$ with a linear model of weights $w$ can be done with the predictor $f_w(x) \triangleq \text{sign}(s(x, w))$. In this situation, a metric of interest quantifying the quality of the classification is given by the margin $m(x, y, w)$, and can be used with the following loss functions:

| Name | Zero-one loss | Hinge loss | Logistic loss |
|---|---|---|---|
| $\text{Loss}(x, y, w)$ | $\mathbb{1}_{\{m(x,y,w) \leqslant 0\}}$ | $\max(1 - m(x, y, w), 0)$ | $\log(1 + e^{-m(x,y,w)})$ |
| Illustration |  |  |  |

❏ **Regression case** — The prediction of a sample $x$ of true label $y \in \mathbb{R}$ with a linear model of weights $w$ can be done with the predictor $f_w(x) \triangleq s(x, w)$. In this situation, a metric of interest quantifying the quality of the regression is given by the margin $\text{res}(x, y, w)$ and can be used with the following loss functions:

| Name | Squared loss | Absolute deviation loss |
|---|---|---|
| $\text{Loss}(x, y, w)$ | $(\text{res}(x, y, w))^2$ | $|\text{res}(x, y, w)|$ |
| Illustration |  |  |

❏ **Loss minimization framework** — In order to train a model, we want to minimize the training loss defined as follows:

$$\text{TrainLoss}(w) = \frac{1}{|\mathcal{D}_{\text{train}}|} \sum_{(x,y) \in \mathcal{D}_{\text{train}}} \text{Loss}(x, y, w)$$

**(https://stanford.edu/~shervine/teaching/cs-221/cheatsheet-reflex-models#non-linear-predictors)**
## Non-linear predictors

❏ **$k$-nearest neighbors** — The $k$-nearest neighbors algorithm, commonly known as $k$-NN, is a non-parametric approach where the response of a data point is determined by the nature of its $k$ neighbors from the training set. It can be used in both classification and regression settings.



| $k = 1$ | $k = 3$ | $k = 11$ |

*Remark: the higher the parameter $k$, the higher the bias, and the lower the parameter $k$, the higher the variance.*

❏ **Neural networks** — Neural networks are a class of models that are built with layers. Commonly used types of neural networks include convolutional and recurrent neural networks. The vocabulary around neural networks architectures is described in the figure below:



Input layer     Hidden layer 1     ...     Hidden layer $k$     Output layer

By noting $i$ the $i^{th}$ layer of the network and $j$ the $j^{th}$ hidden unit of the layer, we have:

$$z_j^{[i]} = {w_j^{[i]}}^T x + b_j^{[i]}$$

where we note $w$, $b$, $x$, $z$ the weight, bias, input and non-activated output of the neuron respectively.

For a more detailed overview of the concepts above, check out the **Supervised Learning cheatsheets (teaching/cs-229/cheatsheet-supervised-learning)**!

**(https://stanford.edu/~shervine/teaching/cs-221/cheatsheet-reflex-models#stochastic-gradient-descent)**
## Stochastic gradient descent

❏ **Gradient descent** — By noting $\eta \in \mathbb{R}$ the learning rate (also called step size), the update rule for gradient descent is expressed with the learning rate and the loss function $\text{Loss}(x, y, w)$ as follows:

$$w \longleftarrow w - \eta \nabla_w \text{Loss}(x, y, w)$$

❐ **Stochastic updates** — Stochastic gradient descent (SGD) updates the parameters of the model one training example $(\phi(x), y) \in \mathcal{D}_{\text{train}}$ at a time. This method leads to sometimes noisy, but fast updates.

❐ **Batch updates** — Batch gradient descent (BGD) updates the parameters of the model one batch of examples (e.g. the entire training set) at a time. This method computes stable update directions, at a greater computational cost.

## (https://stanford.edu/~shervine/teaching/cs-221/cheatsheet-reflex-models#fine-tuning-models)
### Fine-tuning models

❐ **Hypothesis class** — A hypothesis class $\mathcal{F}$ is the set of possible predictors with a fixed $\phi(x)$ and varying $w$:

$$\mathcal{F} = \{f_w : w \in \mathbb{R}^d\}$$

❐ **Logistic function** — The logistic function $\sigma$, also called the sigmoid function, is defined as:

$$\forall z \in ]-\infty, +\infty[, \quad \sigma(z) = \frac{1}{1+e^{-z}}$$

*Remark: we have $\sigma'(z) = \sigma(z)(1 - \sigma(z))$.*

❐ **Backpropagation** — The forward pass is done through $f_i$, which is the value for the subexpression rooted at $i$, while the backward pass is done through $g_i = \frac{\partial \text{out}}{\partial f_i}$ and represents how $f_i$ influences the output.



❐ **Approximation and estimation error** — The approximation error $\epsilon_{\text{approx}}$ represents how far the entire hypothesis class $\mathcal{F}$ is from the target predictor $g^*$, while the estimation error $\epsilon_{\text{est}}$ quantifies how good the predictor $\hat{f}$ is with respect to the best predictor $f^*$ of the hypothesis class $\mathcal{F}$.



❐ **Regularization** — Regularization aims to keep the model from overfitting the data and thus deals with high variance issues. The following table sums up the different types of commonly used regularization techniques:

| LASSO | Ridge | Elastic Net |
|---|---|---|
| • Shrinks coefficients to 0<br>• Good for variable selection | Makes coefficients smaller | Tradeoff between variable selection and small coefficients |
|  |  |  |
| $\ldots + \lambda\|w\|_1$<br>$\lambda \in \mathbb{R}$ | $\ldots + \lambda\|w\|_2^2$<br>$\lambda \in \mathbb{R}$ | $\ldots + \lambda\left[(1-\alpha)\|w\|_1 + \alpha\|w\|_2^2\right]$<br>$\lambda \in \mathbb{R}, \alpha \in [0,1]$ |

❐ **Hyperparameters** — Hyperparameters are the properties of the learning algorithm, and include architecture-related features, the regularization parameter $\lambda$, number of iterations $T$, step size $\eta$, etc.

❐ **Sets vocabulary** — When selecting a model, we distinguish 3 different parts of the data that we have as follows:

| Training set | Validation set | Testing set |
|---|---|---|
| • Model is trained<br>• Usually 80% of the dataset | • Model is assessed<br>• Usually 20% of the dataset<br>• Also called hold-out or development set | • Model gives predictions<br>• Unseen data |

Once the model has been chosen, it is trained on the entire dataset and tested on the unseen test set. These are represented in the figure below:



For a more detailed overview of the concepts above, check out the **Machine Learning tips and tricks cheatsheets (teaching/cs-229/cheatsheet-machine-learning-tips-and-tricks)**

## (https://stanford.edu/~shervine/teaching/cs-221/cheatsheet-reflex-models#unsupervised-learning)
### Unsupervised Learning

The class of unsupervised learning methods aims at discovering the structure of the data, which may have rich latent structures.
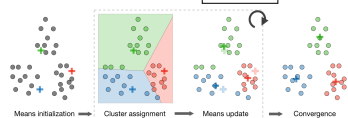
### $k$-means

❐ **Clustering** — Given a training set of n input points $\mathcal{D}_{\text{train}}$, the goal of a clustering algorithm is to assign each point $\phi(x_i)$ to a cluster $z_i \in \{1, ..., k\}$.

❐ **Objective function** — The loss function for one of the main clustering algorithms, $k$-means, is given by:

$$\text{Loss}_{k\text{-means}}(x, \mu) = \sum_{i=1}^{n} \|\phi(x_i) - \mu_{z_i}\|^2$$

❐ **Algorithm** — After randomly initializing the cluster centroids $\mu_1, \mu_2, ..., \mu_k \in \mathbb{R}^d$, the $k$-means algorithm repeats the following step until convergence:

$$z_i = \arg\min_j \|\phi(x_i) - \mu_j\|^2 \quad \text{and} \quad \mu_j = \frac{\sum_{i=1}^{n} 1_{\{z_i = j\}}\phi(x_i)}{\sum_{i=1}^{n} 1_{\{z_i = j\}}}$$



Means initialization ➡ Cluster assignment ➡ Means update ➡ Convergence

### Principal Component Analysis

❐ **Eigenvalue, eigenvector** — Given a matrix $A \in \mathbb{R}^{d \times d}$, $\lambda$ is said to be an eigenvalue of $A$ if there exists a vector $z \in \mathbb{R}^d \setminus \{0\}$, called eigenvector, such that we have:

$$Az = \lambda z$$

❐ **Spectral theorem** — Let $A \in \mathbb{R}^{d \times d}$. If $A$ is symmetric, then $A$ is diagonalizable by a real orthogonal matrix $U \in \mathbb{R}^{d \times d}$. By noting $\Lambda = \text{diag}(\lambda_1, ..., \lambda_d)$, we have:

$$\exists \Lambda \text{ diagonal}, \quad A = U\Lambda U^T$$

*Remark: the eigenvector associated with the largest eigenvalue is called principal eigenvector of matrix $A$.*

❐ **Algorithm** — The Principal Component Analysis (PCA) procedure is a dimension reduction technique that projects the data on $k$ dimensions by maximizing the variance of the data as follows:
- Step 1: Normalize the data to have a mean of 0 and standard deviation of 1.

$$\phi_j(x_i) \leftarrow \frac{\phi_j(x_i) - \mu_j}{\sigma_j} \quad \text{where} \quad \mu_j = \frac{1}{n}\sum_{i=1}^{n} \phi_j(x_i) \quad \text{and} \quad \sigma_j^2 = \frac{1}{n}\sum_{i=1}^{n} (\phi_j(x_i) - \mu_j)^2$$

- Step 2: Compute $\Sigma = \frac{1}{n}\sum_{i=1}^{n} \phi(x_i)\phi(x_i)^T \in \mathbb{R}^{d \times d}$ which is symmetric with real eigenvalues.
- Step 3: Compute $u_1, ..., u_k \in \mathbb{R}^d$ the $k$ orthogonal principal eigenvectors of $\Sigma$, i.e. the orthogonal eigenvectors of the $k$ largest eigenvalues.
- Step 4: Project the data on $\text{span}_{\mathbb{R}}(u_1, ..., u_k)$.

This procedure maximizes the variance among all $k$-dimensional spaces.

Data in feature space ⟶ Find principal components ⟶ Data in principal components space

For a more detailed overview of the concepts above, check out the **Unsupervised Learning cheatsheets (teaching/cs-229/cheatsheet-unsupervised-learning)**!