

Mathematical Notations for Learning from Data Course

Islem Rekik*

Istanbul Technical University

Abstract. Inspiring quotes. *“You’ll need to set small, specific goals to master a skill, but first you’ll want to be sure of the basics.”* Source: *Learn better* by Ulrich Boser.

“Learning is an iterative process that requires that you revisit what you have learnt.” Source: *Make it stick* by Henry L. Roediger III and Mark A. McDaniel.

“The good news is that we now know of simple and practical strategies that anybody can use; at any point in life, to learn better and remember longer: various forms of retrieval practice, such as low-stakes quizzing and self-testing, spacing out practice, interleaving the practice of different but related topics or skills, trying to solve a problem before being taught the solution, distilling the underlying principles or rules that differentiate types of problems, and so on.” Source: *Make it stick* by Henry L. Roediger III and Mark A. McDaniel.

Pre-requisites: Linear algebra. To refresh your memory, you can check 3blue1brown YouTube playlist^a.

Machine Learning Blinks. Check the YouTube link below to watch the lectures^b.

* Corresponding author: irekik@itu.edu.tr; <http://basira-lab.com>

^a https://www.youtube.com/watch?v=fNk_zzaMoSs&list=PLZHQ0b0WTQDPD3MizzM2xVFitgF8hE_ab

^b <https://www.youtube.com/watch?v=HyWmnlahXAA&list=PLug43ldmRSolLDlvQOPzgoJ6wKnfmzimQ>

Table 1: *Major mathematical notations used in lecture 1.*

Mathematical notation	Definition
\mathcal{D}	dataset
n	number of samples in a dataset \mathcal{D}
d	number of features
$\mathbf{x} \in \mathbb{R}^{d \times 1}$	feature vector or data point (sample)
$\mathbf{x}_{feature}^{sample}$	—
$\mathbf{x}^i \in \mathbb{R}^{d \times 1}$	i^{th} sample in the population
$\mathbf{x}_j^i \in \mathbb{R}$	j^{th} feature of i^{th} sample in the population
$\mathcal{D} = \{\mathbf{x}^i, y^i\}_{i=1}^n$	training dataset where $\mathbf{x}^i \in \mathbb{R}^d$ denotes the feature vector for the i^{th} sample and $y^i \in \mathbb{R}$ denotes its score
$\mathbf{X} \in \mathbb{R}^{d \times n}$	data matrix stacking all samples vertically
f	mapping or transformation function to learn
$f : \mathbb{R} \rightarrow \mathbb{R}$	one-to-one mapping
$f : \mathbb{R} \rightarrow \mathbb{R}^n$	one-to-many mapping
$f : \mathbb{R}^p \rightarrow \mathbb{R}$	many-to-one mapping
$f : \mathbb{R}^p \rightarrow \mathbb{R}^m$	many-to-many mapping

Table 2: Major mathematical notations used in lecture 3.

Mathematical notation	Definition
\mathcal{D}	dataset
n	number of samples in a dataset \mathcal{D}
d	number of features
$\mathbf{x} \in \mathbb{R}^{d \times 1}$	feature vector or data point (sample)
$\mathbf{x}_{sample}^{feature}$	–
$\mathbf{x}^i \in \mathbb{R}^{d \times 1}$	i^{th} sample in the population
$\mathbf{x}_j^i \in \mathbb{R}$	j^{th} feature of i^{th} sample in the population
$\Sigma \in \mathbb{R}^{d \times d}$	covariance matrix of data population $\{\mathbf{x}^i\}_{i=1}^n$
$ \mathbf{A} \in \mathbb{R}$	determinant of matrix A
$p(x) = \frac{1}{\sqrt{2\pi}\sigma} \exp(-\frac{1}{2} \frac{\ x-\mu\ _2^2}{\sigma^2})$	probability density function of a variable $x \in \mathbb{R}$
$p(\mathbf{x}) = \frac{1}{(2\pi)^{d/2} \Sigma ^{1/2}} \exp[-\frac{1}{2} (\mathbf{x} - \mu)^T \Sigma^{-1} (\mathbf{x} - \mu)]$	probability density function of a multidimensional variable $\mathbf{x} \in \mathbb{R}^{d \times 1}$
$\mu \in \mathbb{R}^{d \times 1}$	sample mean $\mu = \frac{1}{n} \sum_{i=1}^n \mathbf{x}^i$
$\ \mathbf{x} - \mu\ _{\Sigma^{-1}} = (\mathbf{x} - \mu)^T \Sigma^{-1} (\mathbf{x} - \mu) \in \mathbb{R}$	Mahalanobis distance between \mathbf{x} and μ
$\mathbf{I}_{d \times d} \in \mathbb{R}^{d \times d}$	identity matrix of size $d \times d$
$\ \mathbf{x} - \mu\ _{\mathbf{I}_{d \times d}} = (\mathbf{x} - \mu)^T (\mathbf{x} - \mu) \in \mathbb{R}$	Euclidean distance between \mathbf{x} and μ
$g_i(\mathbf{x}) = \mathbf{w}_i^T \mathbf{x} + \mathbf{w}_{i0}$	discriminant Bayes function for class i when $\Sigma_i = \sigma_i^2 \mathbf{I}$ <u>general case</u> : when $\sigma_i^2 \neq \sigma_j^2$ for classes i and j (i.e., different means $\mu_i \neq \mu_j$ but constant variance for all data features in each class) <u>special case</u> : when $\sigma_i^2 = \sigma_j^2$ for classes i and j (i.e., different means $\mu_i \neq \mu_j$ but constant variances across all classes) (i.e., lines connecting means of different classes are perpendicular to decision boundaries) if $\ln(p(c_i)) = \ln(p(c_j))$, $g_i(\mathbf{x}) = -\ \mathbf{x} - \mu_i\ _2^2$
$g_i(\mathbf{x}) = \mathbf{w}_i^T \mathbf{x} + \mathbf{w}_{i0}$	discriminant Bayes function for class i when $\Sigma_i = \Sigma_j = \Sigma$ (i.e., constant data feature covariance Σ across classes) (i.e., lines connecting means of different classes are not perpendicular to decision boundaries) if $\ln(p(c_i)) = \ln(p(c_j))$, $g_i(\mathbf{x}) = -\frac{1}{2} \ \mathbf{x} - \mu_i\ _{\Sigma^{-1}}^2$
$g_i(\mathbf{x}) = \mathbf{x}^T \mathbf{W} \mathbf{x} + \mathbf{w}_i^T \mathbf{x} + \mathbf{w}_{i0}$	quadratic discriminant function (decision boundaries are nonlinear)

Table 3: Major mathematical notations used in lectures 4.

Mathematical notation	Definition
$\arg \min_{\mathbf{w}} \mathcal{L}(\mathbf{w}) = \frac{1}{n} \sum_{(\mathbf{x}^i, y^i) \in \mathcal{D}} \textcolor{brown}{E}(f_{\mathbf{w}}(\mathbf{x}^i), y^i) + \textcolor{brown}{R}(\dots)$	supervised learning energy cost (loss function)
$\textcolor{brown}{f}$	the mapping function to learn from $\mathbf{x}^i \rightarrow y^i$
$\textcolor{brown}{E}$	the error function between the predicted target by $\textcolor{brown}{f}$ and the ground truth observation y^i
$\textcolor{brown}{R}$	regularization term to avoid overfitting and control model complexity
\mathbf{w}	optimization parameters (weight vector)
$h(w) = l(v) + l'(v)(w - v)$	set of parameters that minimize the loss function $\mathcal{L}(\mathbf{w})$
$l'(w)$	first-order Taylor approximation of the loss function l at point w (in 1-dimensional space)
$h(w) = l(v) + l'(v)(w - v) + \frac{1}{2}l''(v)(w - v)^2$	first derivative of function l evaluated at point w
$l''(w)$	second-order Taylor approximation of the loss function l at point w (in 1-dimensional space)
$h(\mathbf{w}) = l(\mathbf{v}) + \nabla l(\mathbf{w})(\mathbf{w} - \mathbf{v})$	second derivative of function l evaluated at point w
$\mathbf{w} = [w_1 \ w_2 \ \dots \ w_N] \in \mathbb{R}^d$	first-order Taylor approximation of high-dimensional loss function l at vector point $\mathbf{w} \in \mathbb{R}^d$
$\nabla l(\mathbf{w}) \in \mathbb{R}^{d \times 1}$	weight vector to learn
$h(\mathbf{w}) = l(\mathbf{v}) + \nabla l(\mathbf{w})(\mathbf{w} - \mathbf{v}) + \frac{1}{2}(\mathbf{w} - \mathbf{v})^T \nabla^2 l(\mathbf{w})(\mathbf{w} - \mathbf{v})$	gradient vector of the multivariate loss function l at location \mathbf{w}
$\nabla^2 l(\mathbf{w}) \in \mathbb{R}^{d \times d}$	note that $\nabla l(\mathbf{w})^T$ is $\in \mathbb{R}^{1 \times d}$ (row vector)
$l'(w) = 0$	$\nabla l(\mathbf{v}) = [\frac{\partial}{\partial w_1} l(\mathbf{v}) \ \frac{\partial}{\partial w_2} l(\mathbf{v}) \ \dots \ \frac{\partial}{\partial w_d} l(\mathbf{v})]^T$
$\nabla l(\mathbf{w}) = \mathbf{0}_{d \times 1}$	second-order Taylor approximation of high-dimensional loss function l at vector point $\mathbf{w} \in \mathbb{R}^d$
$\mathbf{Q} = \mathbf{Q}^T$	Hessian symmetric matrix of second derivatives of l along all its dimensions (variables)
$\frac{\partial \mathbf{a}^T \mathbf{x}}{\partial \mathbf{x}} = \mathbf{a}$	l is many times differentiable at the vector valued input \mathbf{w}
$\frac{\partial \mathbf{a}^T \mathbf{X} \mathbf{a}}{\partial \mathbf{a}} = 2\mathbf{X} \mathbf{a}$ (for $\mathbf{X} = \mathbf{X}^T$)	stationary point (min, max or saddle) for a 1-dimensional function
$l''(w) > 0$	stationary point (min, max or saddle) for an N -dimensional function
$l''(w) < 0$	(all elements of the gradient are zero)
$\frac{\partial \mathbf{A} \mathbf{x}}{\partial \mathbf{x}} = \frac{1}{2}(\mathbf{A}^T + \mathbf{A})$	\mathbf{Q} is symmetric
$\nabla^2 l(\mathbf{w}) = \frac{1}{2}(\mathbf{Q}^T + \mathbf{Q}) \in \mathbb{R}^{d \times d}$	matrix cookbook ^a (also check ^b)
$\nabla l(\mathbf{w}) = \mathbf{Q} \mathbf{w} + \mathbf{r} \in \mathbb{R}^{d \times 1}$	matrix cookbook
$\mathbf{w}^k = \mathbf{w}^{k-1} - \alpha_k \nabla l(\mathbf{w}^{k-1})$	convex function (facing upward)
	concave function (facing downward)
	matrix cookbook
	is the Hessian matrix of l equal to $l(\mathbf{w}) = \frac{1}{2} \mathbf{w}^T \mathbf{Q} \mathbf{w} + \mathbf{r}^T \mathbf{w} + d$
	$d \in \mathbb{R}$ and $\mathbf{r} \in \mathbb{R}^{d \times 1}$
	is the gradient vector of l equal to $l(\mathbf{w}) = \frac{1}{2} \mathbf{w}^T \mathbf{Q} \mathbf{w} + \mathbf{r}^T \mathbf{w} + d$
	gradient descent for finding the optimal \mathbf{w}

^a <https://www.math.uwaterloo.ca/~hwolkowi/matrixcookbook.pdf>

^b <https://ninova.itu.edu.tr/en/courses/institute-of-science-and-technology/1580/blg-527e/ekkaynaklar?g179937>

Table 4: Major mathematical notations used in lecture 5.

Mathematical notation	Definition
$\mathbf{w}^k = \mathbf{w}^{k-1} - \alpha_k \nabla l(\mathbf{w}^{k-1})$	gradient descent for finding the optimal \mathbf{w}^*
$\nabla^2 l(\mathbf{w}^{k-1}) \mathbf{w}^k = \nabla^2 l(\mathbf{w}^{k-1}) \mathbf{w}^{k-1} - \nabla l(\mathbf{w}^{k-1})$	Newton's method update rule (second order optimization)
$\mathbf{w}^k = \mathbf{w}^{k-1} - [\nabla^2 l(\mathbf{w}^{k-1})]^{-1} \nabla l(\mathbf{w}^{k-1})$	Newton's method update rule when the Hessian matrix $\nabla^2 l(\mathbf{w}^{k-1})$ is invertible

Table 5: Major mathematical notations used in lecture 6 (linear regression and logistic regression).

Mathematical notation	Definition
$\{\mathbf{x}^i, y^i\}_{i=1}^n$	training dataset where $\mathbf{x}^i \in \mathbb{R}^d$ denotes the feature vector for the i^{th} sample and $y^i \in \mathbb{R}$ denotes its score
$\mathbf{x}^T \mathbf{z}$	inner product $\langle \mathbf{x}, \mathbf{z} \rangle$ (or dot product) of vectors $\mathbf{x} \in \mathbb{R}^d$ and $\mathbf{z} \in \mathbb{R}^d$ it outputs a scalar $\in \mathbb{R}$
$\mathbf{x} \mathbf{z}^T$	outer product $\mathbf{x} \otimes \mathbf{z}$ of vectors $\mathbf{x} \in \mathbb{R}^d$ and $\mathbf{z} \in \mathbb{R}^{d'}$ it outputs a matrix of size $d \times d'$
$l(\mathbf{w}, b) = \sum_{i=1}^n (\mathbf{x}^{iT} \mathbf{w} + b - y^i)^2$	least squares loss function for linear regression where $\mathbf{w} \in \mathbb{R}^d, b \in \mathbb{R}$
$l(\mathbf{w}, b) = \sum_{i=1}^n (\tilde{\mathbf{x}}^{iT} \tilde{\mathbf{w}} - y^i)^2$	least squares loss function for linear regression, where $\tilde{\mathbf{x}}^i = [1 \ \mathbf{x}^i]^T$, $\tilde{\mathbf{w}} = [b \ \mathbf{w}]^T$, $\tilde{\mathbf{x}}^i \in \mathbb{R}^{d+1}$, and $\tilde{\mathbf{w}}^* \in \mathbb{R}^{d+1}$
$\nabla l(\tilde{\mathbf{w}}) = \frac{\partial l(\tilde{\mathbf{w}})}{\partial \tilde{\mathbf{w}}} = 0$	setting the gradient to zero to find the optimal solution
$\tilde{\mathbf{w}}^* = [\sum_{i=1}^n \tilde{\mathbf{x}}^i \tilde{\mathbf{x}}^{iT}]^{-1} \sum_{i=1}^n \tilde{\mathbf{x}}^i y^i$	closed form optimum of the least squares regression loss function
$\tilde{\mathbf{w}}^* = [\tilde{\mathbf{X}} \tilde{\mathbf{X}}^T]^{-1} \tilde{\mathbf{X}} \mathbf{y}$	compact closed form where $\tilde{\mathbf{X}}$ vertically stacks all training samples and $\mathbf{y} = [y^1, y^2, \dots, y^n]^T$
$y^i \in \mathbb{R}$	is a scalar denoting the target value (or score) for a single sample i
$\mathbf{y} = [y^1, y^2, \dots, y^n]^T \in \mathbb{R}^n$	is a column vector storing all target values for the n training samples
$\sigma(s) = \frac{1}{1+e^{-s}}$	typical sigmoid function without any scaling and shifting parameters
$\sigma(\mathbf{x}^{iT} \mathbf{w} + b) \sim y^i$	σ maps an input value s to a point $\sigma(s)$ lying on the sigmoid curve
	if the distribution of training samples $\{\mathbf{x}^i\}_{i=1}^n$ slightly vary around the a sigmoid curve
	parameterized by a scaling vector \mathbf{w} and a shifting/translational scalar b
	then identifying the best (\mathbf{w}, b) will identify the best sigmoid curve that fits the data
$l(\mathbf{w}, b) = \sum_{i=1}^n (\sigma(\mathbf{x}^{iT} \mathbf{w} + b) - y^i)^2$	denotes the logistic loss function to minimize
$l(\tilde{\mathbf{w}}) = \sum_{i=1}^n (\sigma(\tilde{\mathbf{x}}^i \tilde{\mathbf{w}}) - y^i)^2$	denotes the compact logistic loss function to minimize
$\frac{\partial \sigma(\tilde{\mathbf{x}}^T \tilde{\mathbf{w}})}{\partial \tilde{\mathbf{w}}} = \sigma(\tilde{\mathbf{x}}^T \tilde{\mathbf{w}}) \times (1 - \sigma(\tilde{\mathbf{x}}^T \tilde{\mathbf{w}})) \times \tilde{\mathbf{x}}^T$	partial derivative of a sigmoid dot product between two vectors $\tilde{\mathbf{x}}$ and $\tilde{\mathbf{w}}$
$\frac{\partial l(\tilde{\mathbf{w}})}{\partial \tilde{\mathbf{w}}} = 0$	we first try to find a closed form of the optimal parameters by setting the gradient of the loss function to 0 to find the optimal parameters $\tilde{\mathbf{w}}^*$
$\frac{\partial l(\tilde{\mathbf{w}})}{\partial \tilde{\mathbf{w}}} = 2 \sum_{i=1}^n (\sigma(\tilde{\mathbf{x}}^{iT} \tilde{\mathbf{w}}) - y^i) \sigma(\tilde{\mathbf{x}}^{iT} \tilde{\mathbf{w}}) (1 - \sigma(\tilde{\mathbf{x}}^{iT} \tilde{\mathbf{w}})) \tilde{\mathbf{x}}^i$	we need to use gradient descent to optimize l since <i>explicitly</i> computing $\tilde{\mathbf{w}}^*$ is not feasible in this case
$l_{reg}(\mathbf{w}, b) = l(\mathbf{w}, b) + \lambda \ \mathbf{w}\ _2^2$	since gradient descent using the original logistic loss function is sensitive to the initial parameters (\mathbf{w}^0, b^0) we introduce a regularization term (squared l_2 norm of the weights to learn if λ is very large, the regularizer drowns out the new loss function $l_{reg}(\mathbf{w}, b)$ λ should be fairly small ~ 0.1 to convexify or smooth the loss function in problematic regions
$\ \mathbf{w}\ _2^2 = \sum_{i=1}^d w_i^2$	the l_2 of a vector is also its Euclidean norm the squared l_2 of a vector is its squared Euclidean norm this is a convex function in \mathbf{w} adding a convex function to a non-convex function convexifies the non-convex function to some extent

Table 6: Major mathematical notations used in lecture 7 (perceptron for classification).

Mathematical notation	Definition
$\{\mathbf{x}^i, y^i\}_{i=1}^n$ $-y^i(\mathbf{x}^{iT}\mathbf{w} + b) < 0$ $-y^i(\mathbf{x}^{iT}\mathbf{w} + b) = 0$ $l_{max}(\mathbf{w}, b) = \sum_{i=1}^n \max(0, -y^i(\mathbf{x}^{iT}\mathbf{w} + b))$ $soft(s_1, s_2) = \log(e^{s_1} + e^{s_2})$ $l_{softmax}(\mathbf{w}, b) = \sum_{i=1}^n soft(0, -y^i(\mathbf{x}^{iT}\mathbf{w} + b))$ $min_{\mathbf{w}, b} l_{softmax}(\mathbf{w}, b)$ $y^{test} = sign(\mathbf{x}^{testT}\mathbf{w}^* + b^*)$ $1 - y^i(\mathbf{x}^{iT}\mathbf{w} + b) \leq 0$ $l_{margin}(\mathbf{w}, b) = \sum_{i=1}^n \max(0, 1 - y^i(\mathbf{x}^{iT}\mathbf{w} + b))$ $l_{softmaxmargin}(\mathbf{w}, b) = \sum_{i=1}^n \log(1 + e^{1-y^i(\mathbf{x}^{iT}\mathbf{w} + b)})$ $lsquaredmargin(\mathbf{w}, b) = \sum_{i=1}^n \max(0, 1 - y^i(\mathbf{x}^{iT}\mathbf{w} + b))^2$	<p>training dataset where $\mathbf{x}^i \in \mathbb{R}^d$ denotes the feature vector for the i^{th} sample and $y^i \in \{-1, 1\}$ denotes its class label</p> <p>criterion for correctly classifying sample \mathbf{x}^i</p> <p>criterion for incorrectly classifying sample \mathbf{x}^i</p> <p>max loss function or perceptron cost for classifying samples $\{\mathbf{x}^i\}_{i=1}^n$</p> <p>issue 1: l_{max} has a trivial solution $\mathbf{w} = 0_{d \times 1}$ and $b = 0$</p> <p>issue 2: l_{max} is non-differentiable at $\mathbf{w} = 0_{d \times 1}$ and $b = 0$</p> <p>soft function approximating $\max(s_1, s_2)$</p> <p>soft max loss function approximating the original perceptron loss</p> <p>$l_{softmax}(\mathbf{w}, b) = \sum_{i=1}^n \log(1 + e^{-y^i(\mathbf{x}^{iT}\mathbf{w} + b)})$</p> <p>solves issue 1: $\mathbf{w} = 0_{d \times 1}$ and $b = 0$ is no trivial solution anymore</p> <p>solves issue 2: is differentiable</p> <p>this minimization problem is also referred to as logistic regression for classification and log-loss support vector machines (see lecture 8)</p> <p>it can be optimized using gradient descent of Newton's method</p> <p>predicting the label y^{test} of a new testing sample \mathbf{x}^{test}</p> <p>checks if sample \mathbf{x}^{test} lies above (+1) or below (-1) the estimated hyperplane parameterized by (\mathbf{w}^*, b^*)</p> <p>\mathbf{w}^* is the normal vector to decision hyperplane</p> <p>margin perceptron (margin of size $2 = 1 + 1$)</p> <p>criterion for correctly classifying sample \mathbf{x}^i constrained by predefined C-sized margin, $C = 2$</p> <p>above the translated hyperplane $\mathbf{x}^T\mathbf{w} + b = 1$</p> <p>or below the translated hyperplane $\mathbf{x}^T\mathbf{w} + b = -1$</p> <p>$\leq$ is used since a data point can lie on either of translated hyperplanes (by +1 or -1) and still get correctly classified</p> <p>margin perceptron loss function for classifying samples $\{\mathbf{x}^i\}_{i=1}^n$</p> <p>issue: non-differentiable</p> <p>soft margin perceptron loss function for classifying samples $\{\mathbf{x}^i\}_{i=1}^n$</p> <p>differentiable and convex \rightarrow we can use gradient descent or Newton's method for optimization</p> <p>squared margin perceptron loss function for classifying samples $\{\mathbf{x}^i\}_{i=1}^n$</p> <p>differentiable and convex \rightarrow we can use gradient descent or Newton's method for optimization</p>

Table 7: *Major mathematical notations used in lecture 8 (support vector machines).*

Mathematical notation	Definition
$\{\mathbf{x}^i, y^i\}_{i=1}^n$	training dataset where $\mathbf{x}^i \in \mathbb{R}^d$ denotes the feature vector for the i^{th} sample and $y^i \in \{-1, 1\}$ denotes its class label
$1 - y^i(\mathbf{x}^{iT} \mathbf{w} + b) \leq 0$	criterion for correctly classifying sample \mathbf{x}^i for margin perceptron defined by hyperplanes $\mathbf{x}^T \mathbf{w} + b = 1$ and $\mathbf{x}^T \mathbf{w} + b = -1$
$\mathcal{L}(\mathbf{w}, b, \alpha)$	Note: both hyperplanes can be translated by $\pm m$, but it is customary to set $m = 1$ $1 - y^i(\mathbf{x}^{iT} \mathbf{w} + b) \leq 0$ is also equivalent to solving $y^i(\mathbf{x}^{iT} \mathbf{w} + b) - 1 \geq 0$ Lagrangian function, $\mathbf{w} \in \mathbb{R}^d, \alpha \in \mathbb{R}^n, \alpha_i \in \mathbb{R}$ for the i^{th} constraint for sample \mathbf{x}^i $\mathcal{L}(\mathbf{w}, b, \alpha) = \frac{1}{2} \ \mathbf{w}\ ^2 - \sum_{i=1}^n \alpha_i (y_i(\mathbf{x}^{iT} \mathbf{w} + b) - 1)$ $\mathcal{L}(\mathbf{w}, b, \alpha) = \frac{1}{2} \mathbf{w}^T \mathbf{w} - \sum_{i=1}^n \alpha_i y_i \mathbf{x}^{iT} \mathbf{w} + b(\sum_{i=1}^n \alpha_i y_i) + \sum_{i=1}^n \alpha_i$
Solving constrained optimization problem	
	Step 1: solve the primal minimization problem: <i>minimize primal variables</i> \mathcal{L} Compute the partial derivatives of \mathcal{L} wrt its primal variables and set them to zero solve $\frac{\partial \mathcal{L}}{\partial \mathbf{w}} = \mathbf{0}_d \rightarrow \mathbf{w}^* = \sum_{i=1}^n \alpha_i y_i \mathbf{x}^i$ (1) solve $\frac{\partial \mathcal{L}}{\partial b} = 0 \rightarrow \sum_{i=1}^n \alpha_i y_i = 0$ (2)
	Step 2: solve the dual maximization problem: <i>maximize Lagrangian multipliers</i> $\mathcal{L}(\mathbf{w}^*, b^*, \alpha)$ Use constraints (1-2) to estimate Lagrangian multipliers $\{\alpha_i\}_{i=1}^n$ by maximizing the dual Lagrangian (compute $\frac{\partial \mathcal{L}}{\partial \alpha} = \mathbf{0}_n$) defines the dual Lagrangian evaluated at optimal parameters \mathbf{w}^*, b^* for SVM to maximize under positivity constraints $\alpha_i \geq 0, i \in 1, \dots, n$ and one equality constraint $\sum_{i=1}^n \alpha_i y_i = 0$ dual Lagrangian to maximize subject to $\alpha_i \geq 0, i \in 1, \dots, n$ and $\sum_{i=1}^n \alpha_i y_i = 0$ $\alpha_i > 0$ for support vectors and $\alpha_i = 0$ otherwise (for non-support vectors)
$\alpha_1^*, \dots, \alpha_n^* = \max_{\alpha_1, \dots, \alpha_n} \mathcal{L}(\mathbf{w}^*, b^*, \alpha)$	
$\mathcal{L}(\mathbf{w}^*, b^*, \alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \mathbf{x}^{iT} \mathbf{x}^j$	
support vectors	

There's another useful learning tool called *self-quizzing*. This learning strategy involves repeatedly recalling and testing yourself on what you've been taught. It's a technique designed to help new ideas stick in your long-term memory.

In fact, research has shown that self-quizzing is 50 percent more effective than some other learning strategies.

A 2006 Washington University study demonstrated just this. Researchers Jeffrey Karpicke and Henry Roediger gave a text to two groups of participants. The first group read it four times. The second group only read it once, but they practiced recalling it three times over.

When Karpicke and Roediger tested them all a few days later, they found that the passage had stuck significantly more in the minds of the group members who'd performed self-quizzing.

Fig. 1: *Blink 1* from 'Learn Better' by Ulrich Boser.