

COMP90051

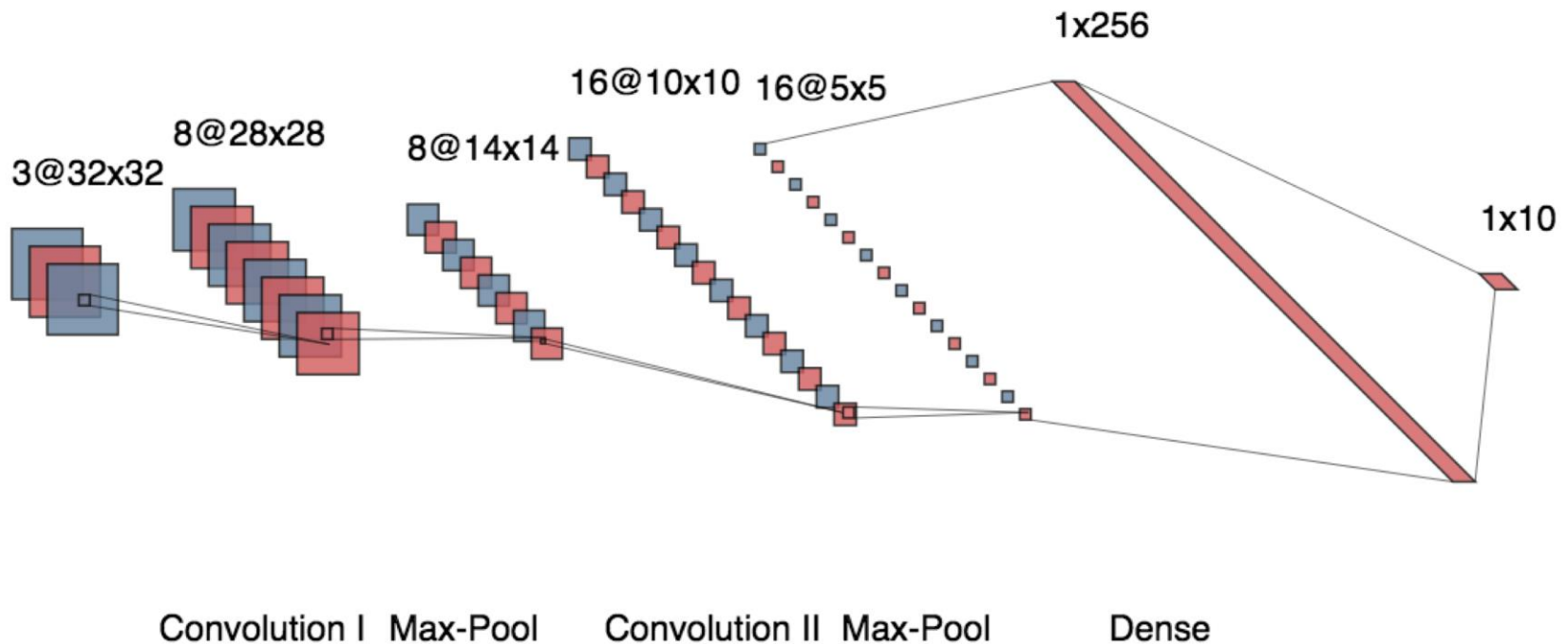
Statistical Machine Learning

Workshop Week 8

Xudong Han

https://github.com/HanXudong/COMP90051_Workshops

Calculate the number of parameters

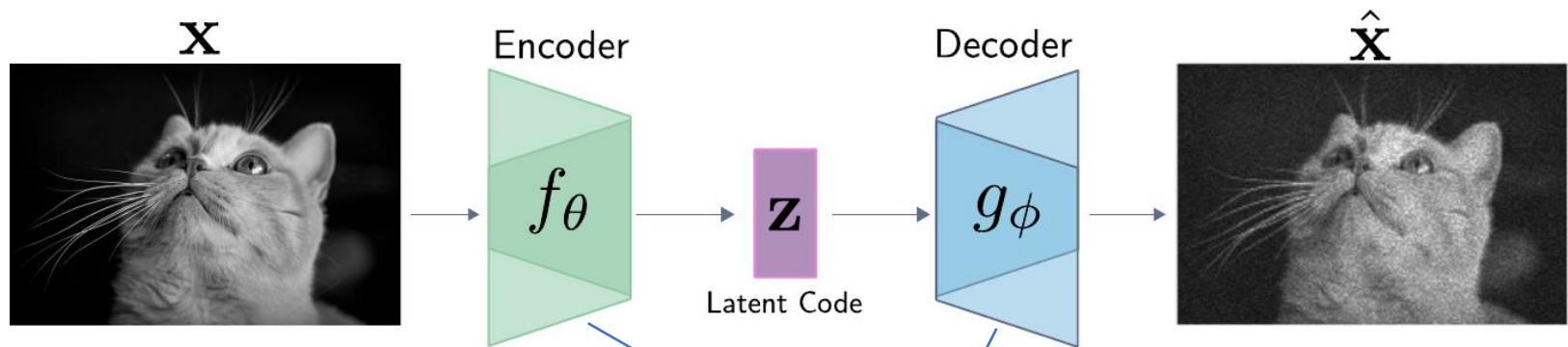


Convolution I: $3 \times 8 \times 5 \times 5 + 8$

<https://pytorch.org/docs/stable/nn.html#convolution-layers>

Dense I: $16 \times 5 \times 5 \times 256 + 256$

Autoencoders



$$\min_{f,g} \sum_k \|\mathbf{x}_k - g \circ f(\mathbf{x}_k)\|^2$$

To do

- SVM hyperparameters: we explore the effect of the penalty parameter and kernel.
- Primal vs. dual: we examine the computational efficiency of the primal and dual formulations in two different scenarios.
- Practice questions.

Kernel Exercises

- Mercer's Theorem
- Positive Semidefinite/ Positive Definite
a symmetric $n \times n$ matrix M is said to be positive semidefinite if for any n nonzero dim vector z , the scalar $z^T M z \geq 0$.
- Eigenvalue
 $Mv = \lambda v$

Eigenvalue

- $Mv = \lambda v$
- Suppose $M = \begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix}$
- How could we compute eigenvalues and eigenvectors?

Positive Semidefinite/Definite

- $[a \quad b] \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} a \\ b \end{bmatrix} = [a + b \quad a + b] \begin{bmatrix} a \\ b \end{bmatrix}$
 $= a^2 + ab + ab + b^2 = (a + b)^2$
- $[a \quad b] \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} a \\ b \end{bmatrix} = a^2 + b^2$
- For a kernel $k(x_i, x_j)$, the full Gram matrix

$$\begin{bmatrix} k(x_1, x_1) & \cdots & k(x_1, x_n) \\ \vdots & \ddots & \vdots \\ k(x_n, x_1) & \cdots & k(x_n, x_n) \end{bmatrix}$$

Mercer's Theorem

– Positive Semidefinite

- Given any $C \in \mathbb{R}^n$
- $$[c_1 \quad \cdots \quad c_n] \begin{bmatrix} k(x_1, x_1) & \cdots & k(x_1, x_n) \\ \vdots & \ddots & \vdots \\ k(x_n, x_1) & \cdots & k(x_n, x_n) \end{bmatrix} \begin{bmatrix} c_1 \\ \vdots \\ c_n \end{bmatrix}$$
- Since we assume that the kernel $k(x_i, x_j)$ is valid kernel

$$C^T K C = \sum_{i,j}^n c_i c_j k(x_i, x_j) \geq 0$$

$$k'(x_i, x_j) = \lambda k(x_i, x_j) \text{ for } \lambda > 0$$

- For a kernel $k'(x_i, x_j)$, the full Gram matrix

$$K' = \begin{bmatrix} \lambda k(x_1, x_1) & \cdots & \lambda k(x_1, x_n) \\ \vdots & \ddots & \vdots \\ \lambda k(x_n, x_1) & \cdots & \lambda k(x_n, x_n) \end{bmatrix}$$

- $C^T K' C = \lambda \sum_{i,j}^n c_i c_j k(x_i, x_j) \geq 0$

- $k'(x_i, x_j)$ is a valid kernel

Extension

- $k'(x_i, x_j) = k_\alpha(x_i, x_j)k_\beta(x_i, x_j)$ is a valid kernel
- can be used for question 2-b

- $k'(x_i, x_j) = \exp(k(x_i, x_j))$

- Taylor series

The Taylor series of a real or complex-valued function $f(x)$ that is infinitely differentiable at a number a is the power series

Kernel Exercises

$$f(x) = \sum_{n=0}^{\infty} \frac{f^{(n)}(a)}{n!} (x - a)^n$$

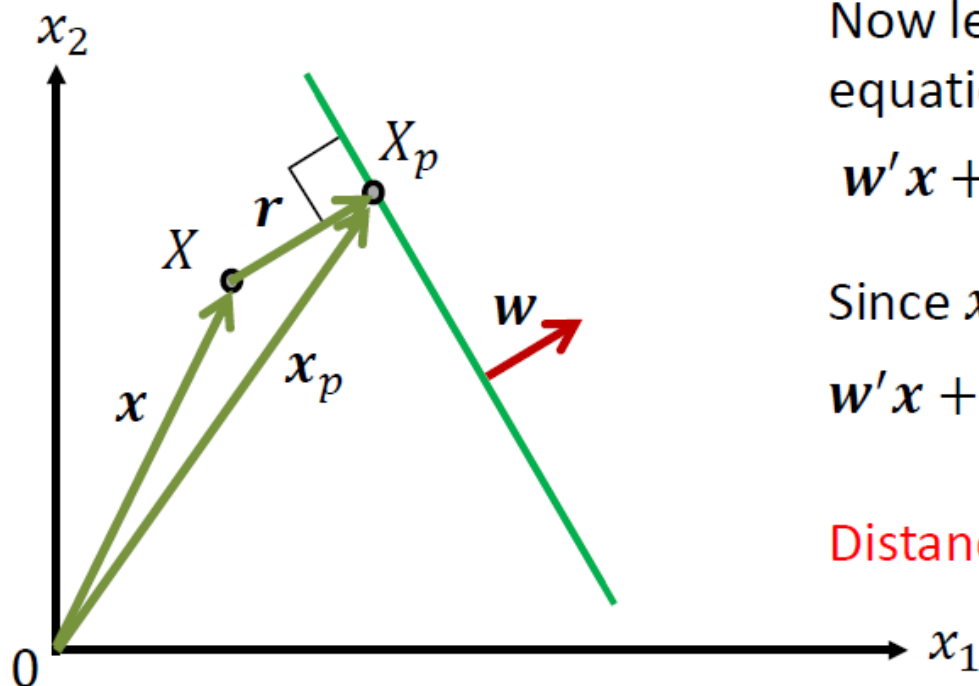
Given $f(x) = \exp(x)$ and $a = 0$, by applying Taylor series:

$$f(x) = \sum_{n=0}^{\infty} \frac{\exp(0)}{n!} (x)^n$$

$$\text{Thus, } \exp(k(x_i, x_j)) = \sum_{n=0}^{\infty} \frac{(k(x_i, x_j))^n}{n!}$$

SVM

- Vectors \mathbf{r} and \mathbf{w} are parallel, but not generally of the same length.
Trivially, $\mathbf{r} = \mathbf{w} \frac{\|\mathbf{r}\|}{\|\mathbf{w}\|}$
- Next, points X and X_p can be viewed as vectors \mathbf{x} and \mathbf{x}_p . By vector addition, we have that $\mathbf{x} + \mathbf{r} = \mathbf{x}_p$ or $\mathbf{x} + \mathbf{w} \frac{\|\mathbf{r}\|}{\|\mathbf{w}\|} = \mathbf{x}_p$



Now let's multiply both sides of this equation by \mathbf{w} and also add b :

$$\mathbf{w}'\mathbf{x} + b + \mathbf{w}'\mathbf{w} \frac{\|\mathbf{r}\|}{\|\mathbf{w}\|} = \mathbf{w}'\mathbf{x}_p + b$$

Since \mathbf{x}_p lies on the boundary, we have

$$\mathbf{w}'\mathbf{x} + b + \|\mathbf{w}\|^2 \frac{\|\mathbf{r}\|}{\|\mathbf{w}\|} = 0$$

Distance is $\|\mathbf{r}\| = -\frac{\mathbf{w}'\mathbf{x} + b}{\|\mathbf{w}\|}$

SVM

- Training data is a collection $\{\mathbf{x}_i, y_i\}$, $i = 1, \dots, n$, where each \mathbf{x}_i is an m -dimensional instance and y_i is the corresponding binary label encoded as -1 or 1
- Given a perfect separation boundary, y_i encode the side of the boundary each \mathbf{x}_i is on
- Thus the distance from the i -th point to a perfect boundary can be encoded as

$$\|\mathbf{r}_i\| = \frac{y_i(\mathbf{w}'\mathbf{x}_i + b)}{\|\mathbf{w}\|}$$

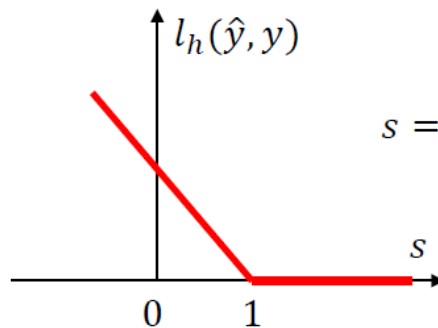
Soft-margin

- Hard-margin SVM loss

$$l_{\infty} = \begin{cases} 0 & 1 - y(\mathbf{w}'\mathbf{x} + b) \leq 0 \\ \infty & \text{otherwise} \end{cases}$$

- Soft-margin SVM loss (**hinge loss**)

$$l_h = \begin{cases} 0 & 1 - y(\mathbf{w}'\mathbf{x} + b) \leq 0 \\ 1 - y(\mathbf{w}'\mathbf{x} + b) & \text{otherwise} \end{cases}$$

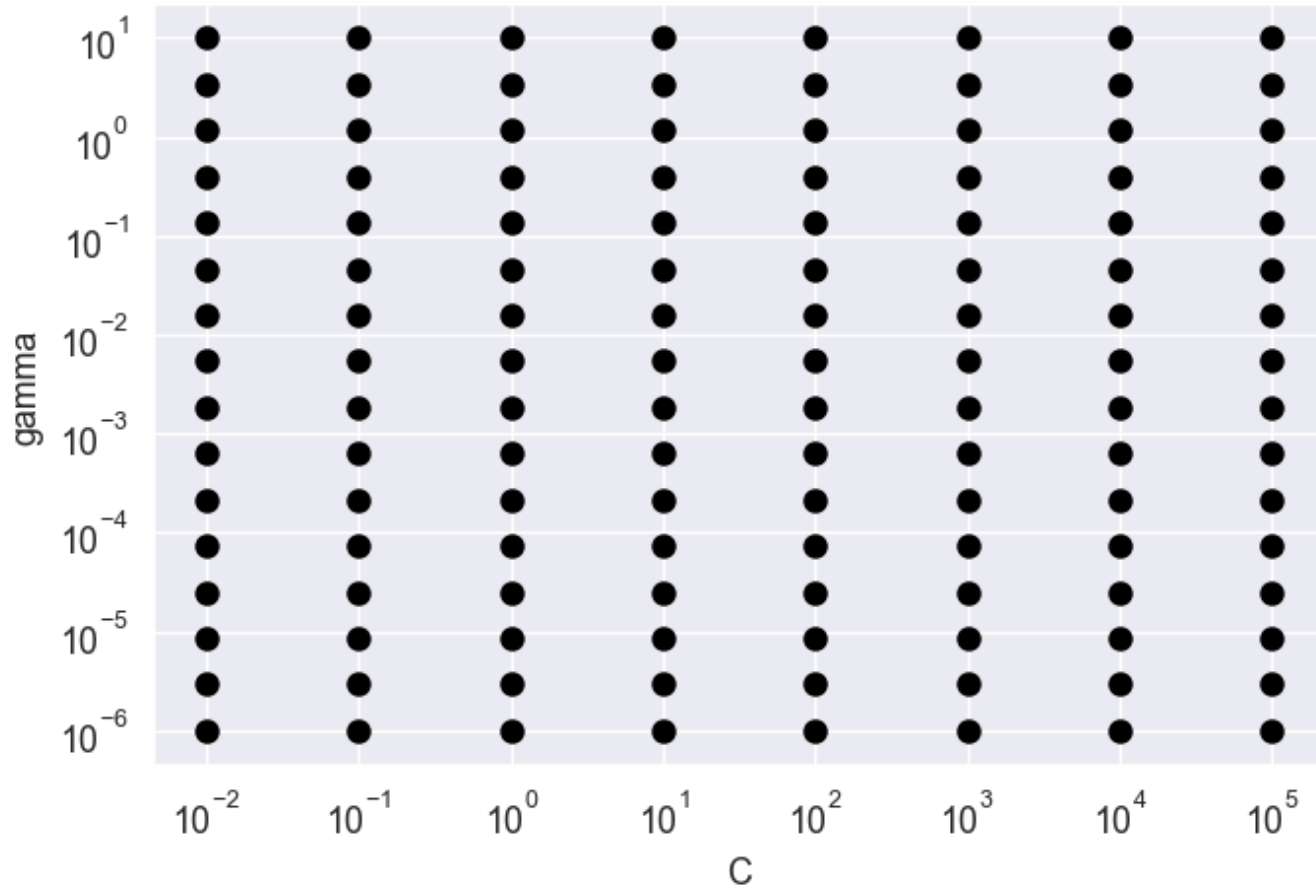


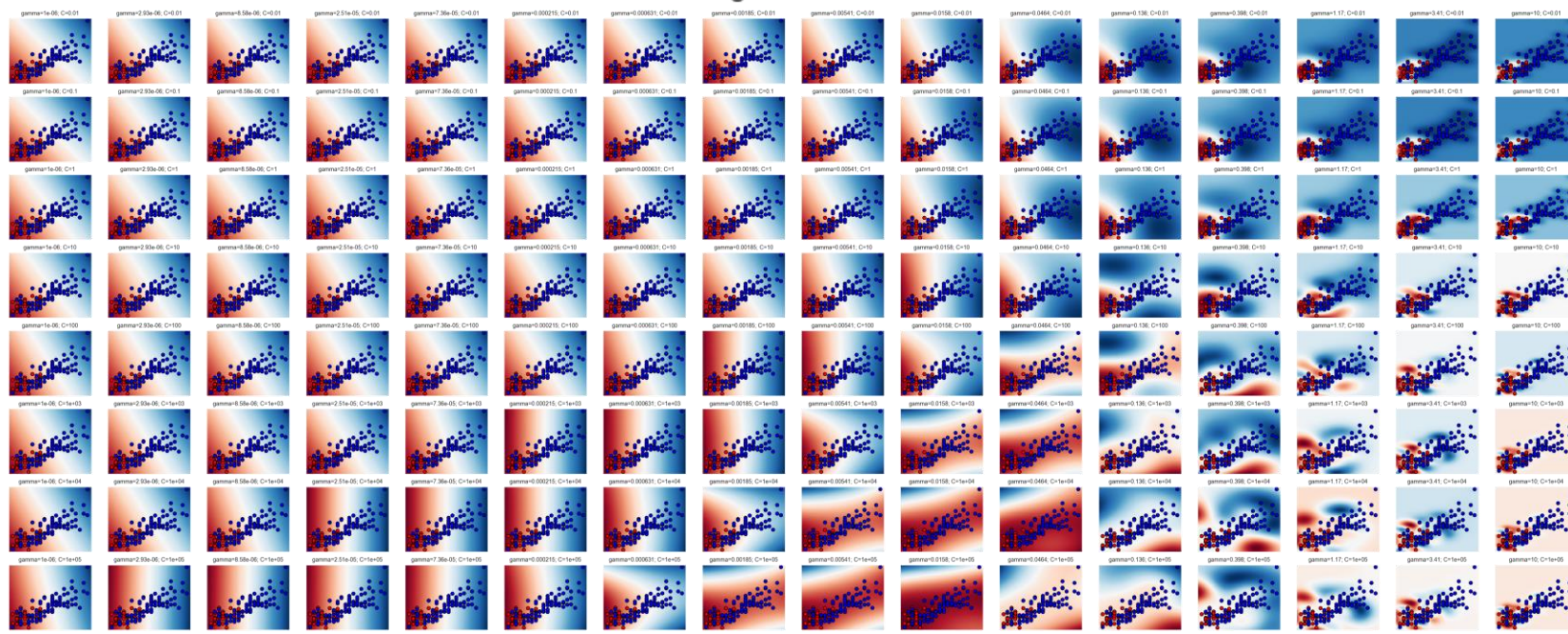
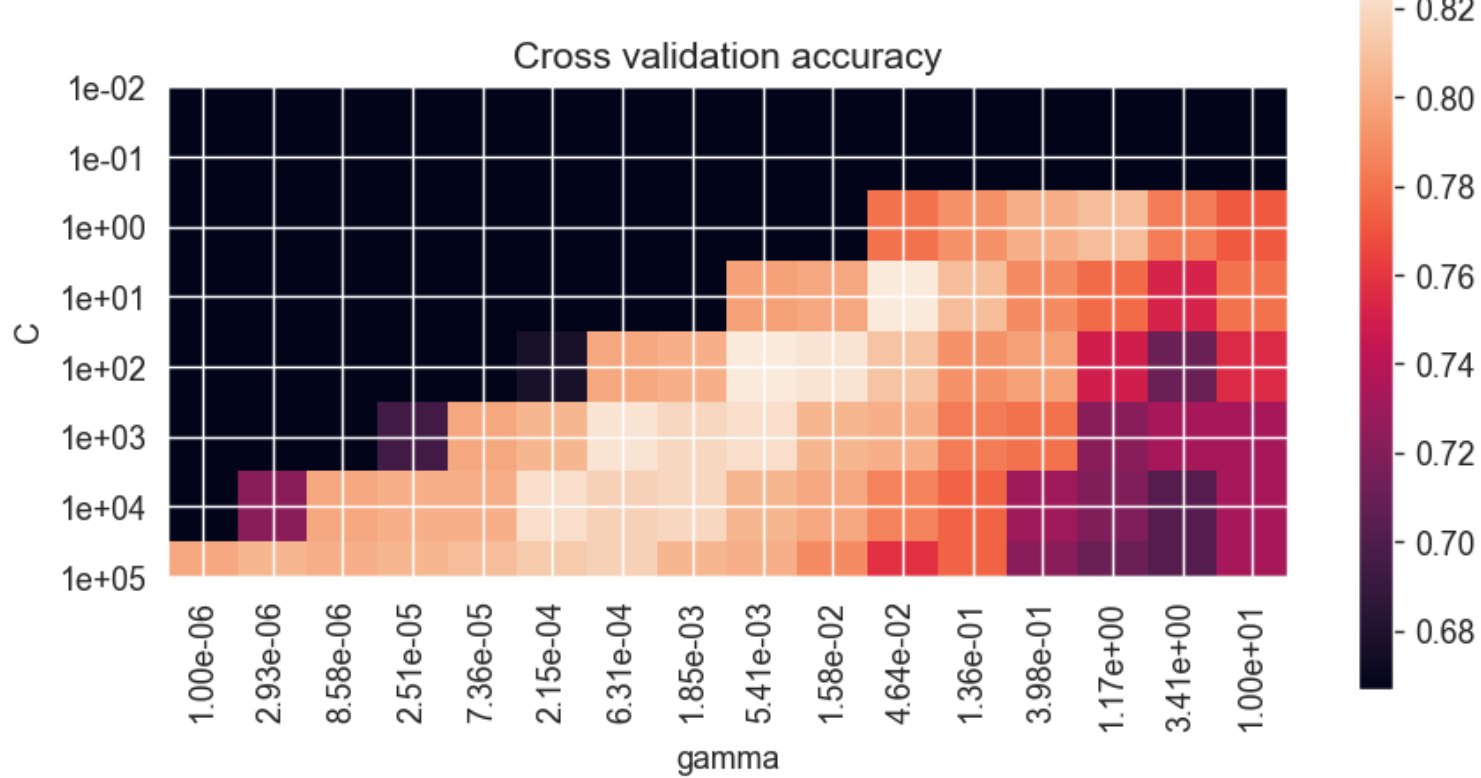
compare this with
perceptron loss

SVC

- <https://scikit-learn.org/stable/modules/svm.html#svc>
- radial basis function (RBF) kernel
- Hyper-parameter
C
gamma

Parameter Grid Search





Primal vs dual

- Introduce auxiliary objective function via auxiliary variables

$$\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}, \mathbf{v}) = f(\mathbf{x}) + \sum_{i=1}^n \lambda_i g_i(\mathbf{x}) + \sum_{j=1}^m v_j h_j(\mathbf{x})$$

Primal constraints became penalties

- * Called the *Lagrangian* function

- * New $\boldsymbol{\lambda}$ and \mathbf{v} are called the *Lagrange multipliers* or *dual variables*

- (Old) **primal program**: $\min_{\mathbf{x}} \max_{\boldsymbol{\lambda} \geq \mathbf{0}, \mathbf{v}} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}, \mathbf{v})$

- (New) **dual program**: $\max_{\boldsymbol{\lambda} \geq \mathbf{0}, \mathbf{v}} \min_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}, \mathbf{v})$

May be easier to solve, advantageous

- Duality theory relates primal/dual:

- * Weak duality: dual optimum \leq primal optimum

- * For convex programs (inc. SVM!) **strong duality**: optima coincide!