**Problem 1 [100 marks].** Consider the $\ell_2$-penalized mean square error (MSE) $J = \frac{1}{2M}\|\mathbf{X}\mathbf{w} - \mathbf{t}\|_2^2 + \lambda\|\mathbf{w}\|_2^2$, where $\mathbf{w} \in \mathbb{R}^{d+1}$ is the weight vector (with an augmented feature) and $M$ is the number of samples.

a) [10 marks] What are the dimensions of $\mathbf{X}$ and $\mathbf{t}$, respectively? *Hint:* No explanation needed.
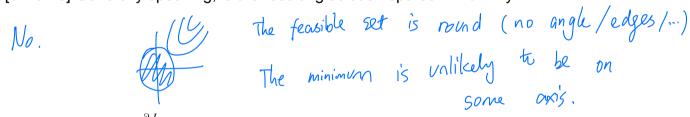
$$X \in \mathbb{R}^{M\times(d+1)} \qquad t \in \mathbb{R}^M$$

b) [10 marks] Transform the soft $\ell_2$-penalty in the form of a hard constraint. Write the formula(s). Draw the feasible set of the constraint (i.e., the region that satisfies the constraint). *Hint:* You may assume $\mathbf{w} \in \mathbb{R}^2$ and use $C$ to represent the constant involved. No proof/explanation needed.

$$\text{minimize} \quad \frac{1}{2M}\|Xw - t\|_2^2$$
$$\text{subject to} \quad \|w\|_2^2 \leq C$$



c) [5 marks] Generally speaking, is the resulting solution sparse? And Why?

No.



The feasible set is round (no angle/edges/...)

The minimum is unlikely to be on some axis.

d) [15 marks] Derive $\frac{\partial J}{\partial \mathbf{w}}$. Please provide derivation steps. *Hints:* You may use either 1) scalar calculus and organize partial derivatives in the vector form, or 2) the matrix calculus identities

- If $\mathbf{A}$ is not a function of $\mathbf{x}$, then $\nabla_\mathbf{x}\mathbf{A}\mathbf{x} = \mathbf{A}^\top$
- If $\mathbf{A}$ is not a function of $\mathbf{x}$ and $\mathbf{A}$ is symmetric, then $\nabla_\mathbf{x}\mathbf{x}^\top\mathbf{A}\mathbf{x} = 2\mathbf{A}\mathbf{x}$

Please also note that $\|\mathbf{x}\|_2^2 = \mathbf{x}^\top\mathbf{x} = \mathbf{x}^\top\mathbf{I}\mathbf{x}$, where $\mathbf{I}$ is an identity matrix.

$$J = \frac{1}{2M}\left(w^TX^T - t^T\right)(Xw - t) + \lambda w^Tw$$

$$= \frac{1}{2M}\left(w^TX^TXw - 2w^TX^Tt + t^Tt\right) + \lambda w^Tw$$

$$\frac{\partial J}{\partial w} = \frac{1}{2M}\left(2X^TXw - 2\cdot X^Tt\right) + 2\lambda w$$

$$= \left(\frac{1}{M}X^TX + 2\lambda I\right)w - \frac{1}{M}X^Tt$$

e) [5 marks] Describe a gradient-based optimization approach in pseudo-code. **Hint:** The gradient has been computed above; no need to repeat. Any variant (e.g., full-batch/mini-batch) is okay.

Loop until convergence:

$$w \leftarrow w - \alpha \cdot \frac{\partial J}{\partial w}$$

f) [10 marks] Alternatively, we may solve the problem by closed-form solution. Compute the closed-form solution for the $\ell_2$-penalized MSE. Provide a few derivation steps.

Set $\frac{\partial J}{\partial w} = 0$

$$\left(\frac{1}{M} X^T X + 2\lambda I\right) w - \frac{1}{M} X^T t = 0$$

$$\left(X^T X + 2\lambda M I\right) w = X^T t$$

$$w = \left(X^T X + 2\lambda M I\right)^{-1} X^T t$$

g) [15 marks] Prove that minimizing $\ell_2$-penalized MSE is equivalent to max *a posterior* (MAP) estimation. **Hint:** Please provide the general principle of MAP estimation (5 marks), appropriate probabilistic assumptions (5 marks), and derivation steps (5 marks).

For a univariate Gaussian distribution, $p(x) \propto \exp\{-\gamma x^2\}$ for some positive $\gamma$. Handling constants is not needed for this question.

General principle    Assume prior $p(w)$, likelihood $p(\mathcal{D}|w)$

$$\text{MAP:} \quad \hat{w}_{MAP} = \underset{w}{\arg\max}\ p(w|\mathcal{D})$$

$$= \underset{w}{\arg\max}\ p(w) \cdot p(\mathcal{D}|w)$$

$$= \underset{w}{\arg\max}\ \log p(w) + \log p(\mathcal{D}|w)$$

For $\ell_2$-penalized MSE    assume

$$p(w_i) \propto \exp\{-\gamma_1 w_i^2\} \qquad \text{for } i=0,\cdots,d$$

$$p(t^{(m)}|x^{(m)}; w) \propto \exp\{-\gamma_2 \left(t^{(m)} - w^T x^{(m)}\right)^2\}$$

$$\text{MAP:} \quad \underset{w}{\text{maximize}}\ \log \prod_{i=0}^{d} \exp\{-\gamma_1 w_i^2\} + \log \prod_{m=1}^{M} \exp\{-\gamma_2(t^{(m)} - w^T x^{(m)})^2\}$$

$$\Longrightarrow \quad \underset{w}{\text{maximize}}\ -\gamma_1 \sum_{i=0}^{d} w_i^2 - \gamma_2 \sum_{m=1}^{M} (t^{(m)} - w^T x^{(m)})^2$$

$$\Longleftrightarrow \quad \underset{w}{\text{minimize}}\ \underbrace{\gamma_1 \|w\|_2^2}_{\ell_2\text{-penalty}} + \underbrace{\gamma_2 \|t^{(m)} - w^T x^{(m)}\|_2^2}_{MSE}$$

h) [5 marks] We know that the mean square error (without $\ell_2$-penalty) yields an unbiased estimate. Why do we prefer $\ell_2$-penalized MSE in some cases? **Hint:** One or a few sentences suffice.

The unbiased estimate assume $t = w^T x + \varepsilon$ for some unknown constant $w$.

But this assumption may not be true.

We need bias-variance tradeoff

i) [10 marks] Fill in the blanks with the word "overfitting" or "underfitting." No explanation needed.

|  | More overfitting or underfitting? |
|---|---|
| Increase $\lambda$ | (a) ___Underfitting___ |
| Exclude some features (i.e., decrease $d$) | (b) ___underfitting___ |

j) [5 marks] Consider labeled data $\mathcal{D}_{\text{label}}$ and test data $\mathcal{D}_{\text{test}}$. Explain why we cannot tune $\lambda$ with $\mathcal{D}_{\text{test}}$. One or a few sentences suffice.

Because $\mathcal{D}_{\text{test}}$ mimicks the deployment, which could only be accessed once in practice.

Tuning $\lambda$ on $\mathcal{D}_{\text{test}}$ leads to over-optimistic performance estimate.

k) [10 marks] Present a correct approach to tune the hyperparameter. Explain how to handle the datasets and provide pseudo code.

Split $\mathcal{D}_{\text{label}} = \mathcal{D}_{\text{train}} \cup \mathcal{D}_{\text{val}}$

For candidate $\lambda$.

  Train $h_\lambda^* = \underset{h \in \mathcal{H}}{\arg\min} \; J_\lambda(h^*, \mathcal{D}_{\text{train}})$

  Validate by $\text{Err}(h_\lambda^*, \mathcal{D}_{\text{val}})$

Pick $\lambda^* = \underset{\lambda}{\arg\min} \; \overline{\text{Err}}(h_\lambda^*, \mathcal{D}_{\text{val}})$

Report test performance $\text{Err}(h_{\lambda^*}^*, \mathcal{D}_{\text{test}})$

**END**

**Scrap paper.** May be detached. Additional paper is available upon request as appropriate.
May be used as answer sheets if you

- Print your name and ID on every answer sheet (including additional sheets) submitted (1 bonus mark)
- Mark your solution and corresponding problem number clearly
- Submit the sheet by the end of the exam