

Optimizers as Penalized Linear Improvement with different norm penalties

In lecture, you saw the locally linear perspective of a neural network and the loss by Taylor expanding the loss around the current value of the parameters. This approximation is only very good in a near neighborhood of those values. One way to proceed with optimization is to consider the size of the neighborhood as a hyperparameter and to bound our update to stay within that neighborhood while minimizing our linear approximation to the loss. You saw in lecture that the choice of norm in defining that neighborhood also matters.

In this problem, you will work out for yourself a slightly different perspective. Instead of treating the norm as a constraint (with the size of the acceptable norm as a hyperparameter), we can do an unconstrained optimization with a weighted penalty that corresponds to the squared norm — where that weight is a hyperparameter.

At each iteration, we wish to maximize linear improvement of the objective (as defined by the dot-product between the gradient and the update) locally regularized by a penalty on the size of the update. This can be expressed (in traditional minimization form) as:

$$1 \\ u = \operatorname{argmin}_{\Delta\theta} g^T \Delta\theta + \frac{\alpha}{2} \|\Delta\theta\|_2^2$$

Linear Improvement

$$+ \\ \alpha \\ d(\Delta\theta)$$

Distance Penalty

$$, (1)$$

where $g = \nabla f(\theta)$ is the gradient of the loss, α is a scalar, and d is a scalar-output distance function

$$d(\Delta\theta) : \mathbb{R}^n \rightarrow \mathbb{R}$$

(a) Let's assume Euclidean distance is the norm that captures our sense of relevant neighborhoods in parameter space. Then our objective can be:

$$u = \operatorname{argmin}_{\Delta\theta} g^T \Delta\theta + \frac{\alpha}{2} \|\Delta\theta\|_2^2$$

2. (2)

What is the analytical solution for u in the above problem? What standard optimizer does this recover?

I've seen quite a few examples of classmates engaging in some light convo with the model while trying to get it to answer the question, so I opted for something a bit different. I proceed with raw and direct prompting, meant to mimic the setting where a student is treating the LLM less as an interactive tutor and more as a tool, a "calculator", for problems.

I paste in the question context and part (a) directly from the pdf of the homework. I make no adjustments to weird formatting and potentially ambiguous notation. I also don't give any context, such as "this is a homework for CS 182 Deep Learning class at Berkeley", and I don't write including any interesting prompting diction that papers like Large Language Models as Optimizers (<https://arxiv.org/abs/2309.03409>) mention, for example "take a deep breath".

🕒 Thought for 69 seconds ▾

- We are given the problem:

$$u = \operatorname{argmin}_{\Delta\theta} \{ \Delta\theta \mid g^T \Delta\theta + (1/\alpha) \|\Delta\theta\|_2^2 \}$$

We can rewrite the objective as:

$$J(\Delta\theta) = g^T \Delta\theta + (1/\alpha) \|\Delta\theta\|_2^2$$

This is a convex function in $\Delta\theta$. To find the minimum, we take the derivative with respect to $\Delta\theta$ and set it to zero.

The gradient of J with respect to $\Delta\theta$ is:

$$\nabla J(\Delta\theta) = g + (2/\alpha) \Delta\theta$$

Setting $\nabla J(\Delta\theta) = 0$:

$$g + (2/\alpha) \Delta\theta = 0$$

$$\Rightarrow (2/\alpha) \Delta\theta = -g$$

$$\Rightarrow \Delta\theta = -(\alpha/2) g$$

Thus, the update rule is: $\Delta\theta = -(\alpha/2) g$.

This is equivalent to gradient descent with a learning rate of $\alpha/2$.

Deepseek nails part (a) quickly and correctly. It rewrites the objective and sets the gradient to zero to find the update rule equation.

However, note that the standard gradient descent update is: $\Delta\theta = -\eta g$, where η is the learning rate.

So, here the learning rate $\eta = \alpha/2$.

But wait, let me double-check the problem. The problem says:

$$u = \operatorname{argmin}_{\{\Delta\theta\}} [g^T \Delta\theta + (1/\alpha) \|\Delta\theta\|_2^2]$$

And we found $\Delta\theta = -(\alpha/2) g$.

But note that sometimes the penalty term is written as $(1/(2\alpha)) \|\Delta\theta\|_2^2$ to make the derivative cleaner. In that case, the objective would be:

$$J(\Delta\theta) = g^T \Delta\theta + (1/(2\alpha)) \|\Delta\theta\|_2^2$$

Then the derivative would be:

$$\nabla J(\Delta\theta) = g + (1/\alpha) \Delta\theta = 0 \Rightarrow \Delta\theta = -\alpha g.$$

So, the standard form often seen is with the $1/2$ factor to cancel the 2 from the derivative of the quadratic term.

In the given problem, the penalty is $(1/\alpha) \|\Delta\theta\|_2^2$, without the $1/2$. So we got $\alpha/2$.

But note: the problem says "recover a standard optimizer". The standard gradient descent update is $\Delta\theta = -\eta g$. So whether it's α or $\alpha/2$, it's still gradient descent, just with a different learning rate.

However, let me check the original problem statement again. It says:

$$u = \operatorname{argmin}_{\{\Delta\theta\}} g^T \Delta\theta + (\alpha) d(\Delta\theta)$$

But then in equation (2), they write:

$$u = \operatorname{argmin}_{\{\Delta\theta\}} g^T \Delta\theta + (1/\alpha) \|\Delta\theta\|_2^2$$

Wait, there's a discrepancy: in the text, they say the penalty is $\alpha d(\Delta\theta)$, but in equation (2), they write $1/\alpha$. Let me read carefully.

The problem says:

"we can do an unconstrained optimization with a weighted penalty that corresponds to the squared norm — where that weight is a hyperparameter."

And then they write the general form (1) as:

$$u = \operatorname{argmin}_{\{\Delta\theta\}} [g^T \Delta\theta + \alpha d(\Delta\theta)]$$

But then for the Euclidean case, they write (2) as:

$$u = \operatorname{argmin}_{\{\Delta\theta\}} [g^T \Delta\theta + (1/\alpha) \|\Delta\theta\|_2^2]$$

So note: in the general form (1), the hyperparameter is α , but in the specific form (2), they use $1/\alpha$. This is just a reparameterization. So in (2), the hyperparameter is $1/\alpha$, meaning that if we let $\lambda = 1/\alpha$, then the penalty is $\lambda \|\Delta\theta\|_2^2$. But then the solution would be $\Delta\theta = -(\lambda/2) g = -(\alpha/2) g$.

So indeed, it's gradient descent with learning rate $\alpha/2$.

But wait, in the lecture, when we use a constraint on the norm, we get a different update (like in projected gradient descent or trust region methods). Here, we are doing an unconstrained optimization with a penalty. And we recover gradient descent.

So the answer is: the analytical solution is $u = -(\alpha/2) g$, and this recovers the standard gradient descent optimizer with learning rate $\alpha/2$.

However, note that in practice, the learning rate is often denoted by η . So we can say that this recovers gradient descent with learning rate $\eta = \alpha/2$.

But let me check the sign: we are minimizing $g^T \Delta\theta + (1/\alpha) \|\Delta\theta\|_2^2$. Since g is the gradient of the loss, minimizing the loss would require moving in the direction opposite to g . So the negative sign is correct.

Alternatively, we can think of it as: we want to minimize the linear approximation of the loss, which is $f(\theta + \Delta\theta) \approx f(\theta) + g^T \Delta\theta$. So we want $g^T \Delta\theta$ to be as negative as possible. But the penalty term prevents $\Delta\theta$ from being too large. So we get a trade-off.

So the solution is indeed $\Delta\theta = -(\alpha/2) g$.

Thus, the standard optimizer is gradient descent.

But wait, there is also the possibility of interpreting this as gradient descent with weight decay? Let me see: the objective is $g^T \Delta\theta + (1/\alpha) \|\Delta\theta\|_2^2$. This is not the same as weight decay because weight decay is applied to the parameters, not the update. This is a penalty on the update size.

So, to summarize:

The analytical solution for u is: $u = -(\alpha/2) g$.

This recovers the standard gradient descent optimizer with learning rate $\alpha/2$.

It also correctly isolates the learning rate as $\alpha/2$ and gets the correct answer. With this info, it correctly concludes that the "recovered standard optimizer" is indeed SGD.

But, probably due to thinking mode being activated, it tries to double check its work even though it didn't run into any problems during its chain of thought up to this point.

Seems to get bogged down on notation here. This is a nice detail that shows Deepseek's robustness to small notation inconsistencies or permutations. In the original question there is no discrepancy, but since I did a rough copy paste, equation (1)'s $1/\alpha$ gets interpreted as α instead. It catches this inconsistency and arrives at the correct conclusion that this is nothing to worry about if you just treat α as an unknown representing the coefficient.

Answer reconfirmed once again.

But note: sometimes the penalty is written as $(1/(2\alpha)) \|\Delta\theta\|_2^2$, which would give $u = -\alpha g$. So the factor of 2 is a matter of convention.

Given the problem as written, we stick with the given penalty term.

So the final answer is:

$$u = -(\alpha/2) \nabla f(\theta)$$

and this is the gradient descent update with learning rate $\alpha/2$.

Solution

Step 1: Rewrite the Objective Function

We want to minimize:

$$J(\Delta\theta) = g^T \Delta\theta + \frac{1}{\alpha} \|\Delta\theta\|_2^2$$

where $g = \nabla f(\theta)$ is the gradient.

After careful thinking, the answer is accurate and formatted clearly and segmented logically.



Step 2: Find the Minimum by Taking the Gradient

Since $J(\Delta\theta)$ is a quadratic function in $\Delta\theta$, we can find its minimum by setting the gradient with respect to $\Delta\theta$ to zero.

$$\nabla_{\Delta\theta} J(\Delta\theta) = g + \frac{2}{\alpha} \Delta\theta$$

Set this equal to zero:

$$g + \frac{2}{\alpha} \Delta\theta = 0$$

Step 3: Solve for $\Delta\theta$

$$\frac{2}{\alpha} \Delta\theta = -g$$

$$\Delta\theta = -\frac{\alpha}{2} g$$

So the optimal update is:

$$u = -\frac{\alpha}{2} g$$

Step 4: Interpret the Result

The update rule:

$$\theta_{\text{new}} = \theta + \Delta\theta = \theta - \frac{\alpha}{2} g$$

is exactly **gradient descent** with learning rate $\eta = \frac{\alpha}{2}$.

Final Answer

$$u = -\frac{\alpha}{2} g$$

This recovers the **standard gradient descent optimizer** with learning rate $\eta = \frac{\alpha}{2}$.

(b) Now, consider an alternative way of capturing local neighborhood size – the squared infinity norm over parameters. Recall that this is defined as $\|x\|_\infty = \max_i |x_i|$. Our objective is now:

$$\begin{aligned} u &= \operatorname{argmin}_{\Delta\theta} \\ &\Delta\theta^T \Delta\theta + \\ &1 \\ &\alpha \|\Delta\theta\|_\infty^2 \\ &\infty, (6) \end{aligned}$$

What is the analytical solution for u in this case? Which optimizer does this correspond to?

Part (b) pasted in here, again with no fillers or reformatting of any kind.

Thought for 301 seconds ✘

- We are given:

$$u = \operatorname{argmin}_{\{\Delta\theta\}} [g^T \Delta\theta + (1/\alpha) \|\Delta\theta\|_\infty^2].$$

But careful: The problem says: "the squared infinity norm" so it's $\|\Delta\theta\|_\infty^2$. But note: $\|\Delta\theta\|_\infty = \max_i |\Delta\theta_{-i}|$. Then squared infinity norm is $(\max_i |\Delta\theta_{-i}|)^2$.

But wait: The problem statement: "the squared infinity norm over parameters. Recall that this is defined as $\|x\|_\infty = \max_i |x_i|$. Our objective is now: $u = \operatorname{argmin}_{\Delta\theta} \Delta\theta^T \Delta\theta + (1/\alpha) \|\Delta\theta\|_\infty^2$." So indeed it's squared infinity norm.

But then we have: $\text{minimize: } g^T \Delta\theta + (1/\alpha) (\max_i |\Delta\theta_{-i}|)^2$.

This is a convex function? But careful: The squared infinity norm is not a quadratic function in the usual sense. But we can solve it.

Let's denote: $M = \|\Delta\theta\|_\infty = \max_i |\Delta\theta_{-i}|$. Then the penalty term is $(1/\alpha) M^2$. And we have constraints: $|\Delta\theta_{-i}| \leq M$ for all i . And we want to minimize $g^T \Delta\theta$ subject to these constraints. But then we also have the penalty on M^2 . So we can optimize over $\Delta\theta$ and M .

So objective: $g^T \Delta\theta + (1/\alpha) M^2$, with constraints: $|\Delta\theta_{-i}| \leq M$ for all i .

We can choose $M \geq 0$. And then for fixed M , the best $\Delta\theta$ to minimize $g^T \Delta\theta$ subject to $|\Delta\theta_{-i}| \leq M$ is: $\Delta\theta_{-i} = -M \operatorname{sign}(g_{-i})$ if $g_{-i} \neq 0$? Actually, careful: $g^T \Delta\theta = \sum_i g_{-i} \Delta\theta_{-i}$. To minimize this given that $\Delta\theta_{-i}$ is constrained to be in $[-M, M]$, the optimal is $\Delta\theta_{-i} = -M \operatorname{sign}(g_{-i})$ if $g_{-i} > 0$ then $\operatorname{sign}(g_{-i})=1$ so $\Delta\theta_{-i} = -M$; if $g_{-i} < 0$ then $\operatorname{sign}(g_{-i})=-1$ so $\Delta\theta_{-i} = M$? But wait, then product: $g_{-i} * (-M \operatorname{sign}(g_{-i})) = -M |g_{-i}|$. So indeed, then $g^T \Delta\theta = -M \sum_i |g_{-i}| = -M \|g\|_1$.

But wait: Is that always optimal? What if some $g_{-i}=0$? Then we can set $\Delta\theta_{-i}=0$. So indeed, the minimal value for fixed M is: $\min_{\{\Delta\theta_{-i}\leq M\}} g^T \Delta\theta = -M \|g\|_1$.

So then the objective becomes: $-M \|g\|_1 + (1/\alpha) M^2$. And we minimize over $M \geq 0$. **That is a quadratic in M : $(1/\alpha) M^2 - \|g\|_1 M$.** The minimum over M (since coefficient positive) is at derivative: $(2/\alpha) M - \|g\|_1 = 0 \Rightarrow M = (\alpha/2) \|g\|_1$. And then the objective value is: $(1/\alpha)(\alpha/2) \|g\|_1^2 - \|g\|_1 (\alpha/2) \|g\|_1 = (\alpha/2)(\alpha^2/4) \|g\|_1^2 - (\alpha/2) \|g\|_1^2 = (\alpha/4) \|g\|_1^2 - (\alpha/2) \|g\|_1^2 = -(\alpha/4) \|g\|_1^2$.

So then the **optimal M is $M^* = (\alpha/2) \|g\|_1$**

And then the optimal $\Delta\theta_{-i}$ are: $\Delta\theta_{-i} = -M^* \operatorname{sign}(g_{-i}) = -(\alpha/2) \|g\|_1 \operatorname{sign}(g_{-i})$. But note: $\operatorname{sign}(g_{-i})$ times $\|g\|_1$? Actually, that gives: $\Delta\theta_{-i} = -(\alpha/2) (\|g\|_1) \operatorname{sign}(g_{-i})$. But $\|g\|_1$ is part of $\|g\|_1$. But note: This update is: $\Delta\theta_{-i} = -(\alpha/2) (\sum_j |g_{-j}|) \operatorname{sign}(g_{-i})$. That is like a normalized sign? Actually, it's similar to the update in coordinate descent? Alternatively, it is like: $\Delta\theta = -(\alpha/2) \|g\|_1 \operatorname{sign}(g)$, where s is the vector of signs of g . **That is not exactly the same as gradient descent because gradient descent would be $\Delta\theta = -\eta g$.** Here it's proportional to the sign of g , but the step size is constant across coordinates? Actually, it is: $\Delta\theta_{-i} = -(\alpha/2) \|g\|_1 \operatorname{sign}(g_{-i})$. So the magnitude of the update in each coordinate is the same: $(\alpha/2) \|g\|_1$. That is like a version of sign descent? But wait, there is an optimizer known as "normalized gradient descent" or something? Alternatively, consider: If we use L^∞ norm constraint instead of squared? But here we have squared L^∞ norm penalty.

Wait, but the problem statement: "Which optimizer does this correspond to?" Possibly it is related to the sign descent optimizer? Actually, there is an optimizer called "**Signed gradient descent**" where the update is proportional to the sign of the gradient. But here the step size is $(\alpha/2) \|g\|_1$. **That is not exactly the same as signSGD?** SignSGD update is: $\Delta\theta = -\eta \operatorname{sign}(g)$. But here the step size is $\eta = (\alpha/2) \|g\|_1$. So it's sign gradient descent but with a step size that is proportional to the L^1 norm of the gradient. **That is like an adaptive step size?**

Actually, wait, there is also the method called "Adaptive Gradient Algorithm" (AdaGrad) but that is different.

Maybe we can reinterpret: $u = \operatorname{argmin}_{\Delta\theta} [g^T \Delta\theta + (1/\alpha) \|\Delta\theta\|_\infty^2]$. The solution we got: $\Delta\theta_{-i} = -(\alpha/2) \|g\|_1 \operatorname{sign}(g_{-i})$. That is equivalent to: $\Delta\theta = -(\alpha/2) \|g\|_1 \operatorname{sign}(g)$. That is like: $\Delta\theta = -(\alpha/2) * (\|g\|_1 / \|g\|_\infty)$. No.

Alternatively, consider if we had used the L^∞ norm (not squared) as penalty? That would

Correctly identifies convexity (sum of linear + convex penalty) but hesitates on "quadratic" due to non-Euclidean norm. Which shows a correct understanding of convexity.

Nice reformulation here. Converts non-smooth objective into constrained scalar optimization via auxiliary variable M . Leverages epigraph form implicitly.

Accurately captures the quadratic in M equation resulting from the above simplifications

Notices the "extra" $\|g\|_1$ term that makes it not standard gradient descent.

Correctly finds the update equation with some thought roaming but not much (i.e., just a bit more than half a page here).

Moving onto the second question of which optimizer it corresponds to. Immediately gets the right intuition for signGD.

Perfectly finds the answer to this question, almost verbatim to the solution key. If we could limit Deepseek's thinking tokens (like Qwen) we would've gotten the correct answer at this point. But... as is the case with Deepseek, it loves to think more and double check itself like ten times.

be: $\mathbf{g}^T \Delta\theta + (1/a) \|\Delta\theta\|_\infty$. Then for fixed M, the linear part gives $-M\|\mathbf{g}\|_1$, and penalty gives $(1/a) M$, so total: $-M\|\mathbf{g}\|_1 + (1/a) M = M(1/a - \|\mathbf{g}\|_1)$. Then if $\|\mathbf{g}\|_1 > 1/a$, then we want M as large as possible? That would be unbounded? So that doesn't work. So squared norm is necessary to get a finite solution.

Wait, maybe there is another interpretation: The problem says "Optimizers as Penalized Linear Improvement with different norm penalties". In part (a) we got gradient descent. In part (b) we get an optimizer that uses the sign of the gradient. There is an optimizer called "Signed gradient descent" or "SignSGD". But SignSGD is usually: $\Delta\theta = -\eta \text{sign}(\mathbf{g})$. Here we have $\Delta\theta = -(\alpha/2) \|\mathbf{g}\|_1 \text{sign}(\mathbf{g})$. So the step size is not constant; it depends on the L1 norm of the gradient. That is actually the same as: $\Delta\theta = -(\alpha/2) * (\|\mathbf{g}\|_1) * \text{sign}(\mathbf{g})$. But note that $\|\mathbf{g}\|_1 \text{sign}(\mathbf{g})$ is not the same as \mathbf{g} . However, if we consider the L1 norm of \mathbf{g} , then the update in each coordinate is proportional to the sign of the gradient, but the magnitude is the same for all coordinates. That is like a normalized update where the step size is determined by the total magnitude of the gradient. That is reminiscent of the update in the "Frank-Wolfe" algorithm? Or maybe it's related to "Coordinate descent" with a specific step? Alternatively, it might be related to "Mirror descent" with the L^∞ norm as the potential? Actually, mirror descent with the L^∞ norm potential leads to updates that are like sign gradients? Let's check: Mirror descent using the entropy function (for simplex) gives multiplicative updates. But here the domain is all R^n . The dual norm of L^∞ is L^1 . There is a connection: The proximal step with L^∞ norm squared? Actually, the proximal operator of the squared L^∞ norm? The function $d(x) = (1/2)\|x\|^2$ is not separable. Its proximal operator is given by something like: $\text{prox}_{\lambda\|\cdot\|^\infty/2}(v) = ?$ That might yield a sign vector. I recall that the proximal operator of the L^∞ norm (not squared) is the soft-thresholding. But here it's squared.

Maybe we can re-derive carefully:

We want to minimize: $F(\Delta\theta) = \mathbf{g}^T \Delta\theta + (1/a) (\max_i |\Delta\theta_i|)^2$.
Let $M = \max_i |\Delta\theta_i|$. Then necessarily, for each i, $|\Delta\theta_i| \leq M$. And the best choice to minimize $\mathbf{g}^T \Delta\theta$ is to set $\Delta\theta_i = -M \text{sign}(\mathbf{g}_i)$ (if $\mathbf{g}_i \neq 0$; if $\mathbf{g}_i=0$, then any value in $[-M, M]$ works, but then we can set 0 to not increase M). Actually, if $\mathbf{g}_i=0$, then that coordinate doesn't affect the linear term. But then we want to minimize M? Actually, M is defined as the maximum of $|\Delta\theta_i|$. So if we set $\Delta\theta_i=0$ for those with $\mathbf{g}_i=0$, then M is determined by the ones with nonzero g. So indeed, we can take $\Delta\theta_i = -M \text{sign}(\mathbf{g}_i)$ for all i with $\mathbf{g}_i \neq 0$, and for $\mathbf{g}_i=0$, we can set $\Delta\theta_i=0$. But then M is the maximum of these $|\Delta\theta_i|$, so M will be M itself because for nonzero, $|\Delta\theta_i| = M$. So then M is indeed the maximum, so it's consistent. So then F becomes: $\sum_i \mathbf{g}_i (-M \text{sign}(\mathbf{g}_i)) + (1/a) M^2 = -M \sum_i |\mathbf{g}_i| + (1/a) M^2 = -M\|\mathbf{g}\|_1 + (1/a) M^2$.

Then minimize over $M \geq 0$. Derivative: $-\|\mathbf{g}\|_1 + (2/a) M = 0 \Rightarrow M = (a/2) \|\mathbf{g}\|_1$.

So then $\Delta\theta_i = -(\alpha/2) \|\mathbf{g}\|_1 \text{sign}(\mathbf{g}_i)$. So that is the solution.

Now, what optimizer is this? It is an update that is proportional to the sign of the gradient. But note that standard sign descent uses a constant step size: $\Delta\theta = -\eta \text{sign}(\mathbf{g})$. Here the step size is $\eta = (\alpha/2) \|\mathbf{g}\|_1$. So it's like an adaptive sign descent where the step size is scaled by the L1 norm of the gradient. That is actually equivalent to: $\Delta\theta = -(\alpha/2) (\|\mathbf{g}\|_1) \text{sign}(\mathbf{g})$. But if we write componentwise: $\Delta\theta_i = -(\alpha/2) (\sum_j |\mathbf{g}_{j,i}|) \text{sign}(\mathbf{g}_i)$. This is not the same as gradient descent. There is an optimizer called "Normalized Sign Descent" or something? Alternatively, consider if we had used the L1 norm penalty? That would give something like: $\Delta\theta = -(\alpha/2) \text{sign}(\mathbf{g})$ as well? Let's check: If we had $d(\Delta\theta) = \|\Delta\theta\|_1$, then the objective would be: $\mathbf{g}^T \Delta\theta + (1/a) \|\Delta\theta\|_1$. That is not quadratic, and its minimization gives soft-thresholding. That is like LASSO. But that's not it.

Wait, maybe there is a connection to the "infinity norm" and the "dual norm" idea. In optimization, when we use a constraint $\|\Delta\theta\|_\infty \leq \epsilon$, then the best linear improvement is achieved by $\Delta\theta = -\epsilon \text{sign}(\mathbf{g})$. So that is exactly the update from the constrained formulation: maximize $\mathbf{g}^T \Delta\theta$ subject to $\|\Delta\theta\|_\infty \leq \epsilon$ gives $\Delta\theta = -\epsilon \text{sign}(\mathbf{g})$. So if we set $\epsilon = (a/2) \|\mathbf{g}\|_1$? That doesn't match because then ϵ would be adaptive. Actually, in the constrained formulation, ϵ is a hyperparameter. Here, in the penalized formulation, the equivalent ϵ becomes $(a/2) \|\mathbf{g}\|_1$. So it's like the step size is determined by the gradient's L1 norm. That is interesting.

Maybe this optimizer is known as "SignSGD with momentum"? Not sure.

I recall that there is an optimizer called "Adam" that uses signs in a way? Actually, Adam uses an estimate of the first moment which is like a sign if the gradient is consistent. But not exactly.

Wait, the problem likely expects: For squared L2 norm penalty, we get gradient descent. For squared L^∞ norm penalty, we get sign gradient descent. But sign gradient descent is usually $\Delta\theta = -\eta \text{sign}(\mathbf{g})$. But here we have an extra factor of $\|\mathbf{g}\|_1$. So if we define $\eta = (\alpha/2) \|\mathbf{g}\|_1$, then it is sign gradient descent with a step size that depends on the L1 norm of the gradient. That is not standard sign descent because standard sign descent uses a constant step size. However, if we consider the normalized sign descent, where we normalize the update by the L1 norm, then actually $\text{sign}(\mathbf{g})$ is the direction of the gradient in the L^∞ norm? Because the dual of L^1 is L^∞ . Actually, wait: The gradient descent direction in the L^∞ norm geometry? There is something called "mirror descent" with the L1 norm potential leads to L^∞ norm constraints. Specifically, if we use the entropy function for the simplex, but for the

It seems like the fact that the answer for part (b) isn't as clean as part (a) is bothering Deepseek. I.e., there is no clean descent name to call it like SGD. The answer being a twist on signSGD (with adaptive step size) is making Deepseek doubt itself repeatedly.

Double checked successfully.

Interesting detail here. Deepseek seems to be trying to meta-game the solution by predicting what the question makers would like to see. I guess this psychological tactic is like that "oh I've only had (a) (b) (d) choices for every question and not a single (c) so surely I must've gotten at least one question wrong" belief in multi choice exams.

whole space, if we use the potential $\Psi(x) = (1/2)\|x\|_1^2$? That doesn't sound standard.

Alternatively, consider the following: The update we got is $\Delta\theta_i = -(\alpha/2)\|g\|_1 \text{sign}(g_i)$. This can be rewritten as $\Delta\theta_i = -(\alpha/2)(\sum_j |g_j|) \text{sign}(g_i)$. Notice that if we consider the vector whose components are all the same magnitude, then this update is proportional to the sign of the gradient. So it is essentially a signed update. And the step size is $(\alpha/2)\|g\|_1$. So it is like: $\Delta\theta = -(\alpha/2)\|g\|_1 \text{sign}(g)$. This is equivalent to: $\Delta\theta = -(\alpha/2) * (\|g\|_1/\|g\|_\infty)$.

I think the intended answer is that this corresponds to the "SignSGD" optimizer. But wait, SignSGD is usually: $\theta_{t+1} = \theta_t - \eta \text{sign}(\nabla f(\theta_t))$. So here η would be $(\alpha/2)\|g\|_1$. But that is not a constant learning rate; it depends on the gradient. However, if we absorb the $\|g\|_1$ into α , then α is a hyperparameter. But then the update becomes: $\Delta\theta = -(\alpha/2)\|g\|_1 \text{sign}(g)$. If we set $\alpha' = \alpha/2$, then it's $-\alpha'\|g\|_1 \text{sign}(g)$. That is not standard because in SignSGD, the learning rate is usually set independently of the gradient magnitude. But maybe in the context of this problem, they consider that the hyperparameter α can be adjusted. Actually, note that in the constrained formulation (which was discussed in lecture), the update for the L_∞ norm constraint was: $\Delta\theta = -\varepsilon \text{sign}(g)$. So that is SignSGD. And here we got a penalty formulation that gives $\Delta\theta = -(\alpha/2)\|g\|_1 \text{sign}(g)$. So these are different unless $\|g\|_1$ is constant? That is not generally true.

Maybe I made a mistake: The objective is: $g^T \Delta\theta + (1/\alpha) \|\Delta\theta\|_\infty^2$. But wait, the norm is squared infinity norm. But careful: $\|\Delta\theta\|_\infty^2$ is not the same as $(\|\Delta\theta\|_\infty)^2$? Actually, it is: $\|\Delta\theta\|_\infty^2 = (\max_i |\Delta\theta_i|)^2$. So that is correct.

What if the penalty was $(1/\alpha) \|\Delta\theta\|_\infty^2$? That is what we used.

Maybe the intended penalty was actually $(1/\alpha) \|\Delta\theta\|_\infty$ (without square)? But then the problem says "squared infinity norm". So it's squared.

Let me check the lecture note idea: "Optimizers as Penalized Linear Improvement with different norm penalties". In lecture, they might have considered: $u = \operatorname{argmin}_{\{\Delta\theta\}} g^T \Delta\theta$ subject to $\|\Delta\theta\| \leq \delta$. And then for L2 norm, that gives gradient descent with step size depending on δ . For L_∞ norm, that gives $\Delta\theta = -\delta \text{sign}(g)$. So that is SignSGD. And here we are doing a penalty version instead of constraint. So for L2 penalty, we got gradient descent. For L_∞ penalty, we should get something that corresponds to the constraint version with a particular δ . But wait, in the penalty method, the hyperparameter α plays the role of the inverse of the penalty strength. In the constraint method, δ is the constraint size. There is a relation between α and δ via Lagrange multipliers. Specifically, if we solve the constrained problem: minimize $g^T \Delta\theta$ subject to $\|\Delta\theta\|_\infty \leq \delta$, the solution is $\Delta\theta = -\delta \text{sign}(g)$.

And if we solve the penalized problem: minimize $g^T \Delta\theta + (1/\alpha) \|\Delta\theta\|_\infty^2$, we got $\Delta\theta = -(\alpha/2)\|g\|_1 \text{sign}(g)$. For these to be equivalent, we would need $\delta = (\alpha/2)\|g\|_1$. But δ is a hyperparameter that is fixed, while $\|g\|_1$ changes with iteration. So they are not the same optimizer. So the penalized formulation with squared L_∞ norm gives an optimizer that adapts the step size proportional to the L1 norm of the gradient. That is like an adaptive sign descent. Is there a name for that? It might be related to "AdaGrad" but AdaGrad uses the square root of the sum of squares of gradients, not the L1 norm.

Wait, maybe I misinterpreted: The problem says: "the squared infinity norm" but then it writes $\|\Delta\theta\|_\infty^2$. That could be interpreted as the square of the infinity norm, which is what I did. But what if it meant the squared norm in the sense of the L_∞ norm squared? That is the same.

Alternatively, could it be that the penalty is $(1/\alpha) \|\Delta\theta\|_\infty^2$, where $\|\Delta\theta\|_\infty$ is defined as the maximum absolute value, but then squaring it gives a function that is not strictly a norm squared because the infinity norm is not Euclidean. But it is still a valid function.

I will re-read the problem statement: "Now, consider an alternative way of capturing local neighborhood size – the squared infinity norm over parameters. Recall that this is defined as $\|x\|_\infty = \max_i |x_i|$. Our objective is now: $u = \operatorname{argmin} \Delta\theta \ g^T \Delta\theta + (1/\alpha) \|\Delta\theta\|_\infty^2$ ". So it is clearly $(\max_i |\Delta\theta_i|)^2$.

Maybe the intended optimizer is "Coordinate Descent" but that updates one coordinate at a time. This updates all coordinates simultaneously with the same magnitude. That is like a version of "Signed Gradient Descent" with a step size that is the same for all coordinates. That is sometimes used in practice for its robustness to gradient scale.

I will check literature: There is an optimizer called "Normalized Gradient Descent" where the update is $\Delta\theta = -\eta g/\|g\|$. That is different. There is also "Gradient Descent with L_∞ -norm trust region" which gives $\Delta\theta = -\eta \text{sign}(g)$. So I lean towards the idea that the intended answer is that this recovers the SignSGD optimizer. But wait, in the constraint formulation, SignSGD comes from L_∞ constraint. And here we have a penalty formulation that is different. But maybe the problem expects: For L2 penalty, gradient descent; for L_∞ penalty, sign gradient descent. Even though the step size here is $(\alpha/2)\|g\|_1$, if we view α as a hyperparameter, then effectively the update is proportional to $\text{sign}(g)$. Because $\|g\|_1$ is a scalar that depends on the gradient. But then the learning rate is not constant; it changes with the gradient's L1 norm. However, in many optimizers, the learning rate is adjusted adaptively. But SignSGD is usually presented with a constant learning rate.

Repeats these lines many times: "oh it should be SignSGD, but the varying gradient term makes it not actually SignSGD, then what is it? Oh might've gotten something wrong. Let's check again...."

Similar repetition seen here....

Maybe there is a mistake: What if the penalty was $(1/\alpha) \|\Delta\theta\|_\infty$ (not squared)? Then the objective would be: $g^T \Delta\theta + (1/\alpha) \|\Delta\theta\|_\infty$. Then for fixed $M = \|\Delta\theta\|_\infty$, the best linear term is $-M\|g\|_1$, so total: $-M\|g\|_1 + (1/\alpha) M = M(1/\alpha - \|g\|_1)$. This is linear in M . If $\|g\|_1 > 1/\alpha$, then the coefficient is negative, so we would want M as large as possible, so the minimum is not attained (goes to $-\infty$). If $\|g\|_1 < 1/\alpha$, then the coefficient is positive, so the minimum is at $M=0$. If $\|g\|_1 = 1/\alpha$, then any M gives the same value. So that doesn't yield a nice optimizer. So the squared norm is necessary to get a finite solution.

Given that part (a) gave gradient descent, and part (b) is likely to give an optimizer that is different from gradient descent, and the most common one associated with the L_∞ norm is sign descent, I will answer that it corresponds to Sign Gradient Descent.

But let me check the form: $\Delta\theta = -(a/2) \|g\|_1 \text{sign}(g)$. If we write this componentwise: $\Delta\theta_i = -(a/2) (\sum_j |g_{ji}|) \text{sign}(g_{ii})$. This is actually equivalent to: $\Delta\theta_i = -(a/2) (\sum_j |g_{ji}|) * (g_{ii}/|g_{ii}|)$ if $g_{ii} \neq 0$. So it's like a per-coordinate learning rate that is $(a/2)(\sum_j |g_{ji}|)/|g_{ii}|$. That is similar to AdaGrad if the historical sum of gradients is replaced by the current L1 norm? Not exactly.

Wait, maybe I should consider the possibility that the penalty term is misinterpreted. Sometimes, the squared norm means the square of the norm, but there is also the notion of "norm squared" which for the infinity norm is not standard because the infinity norm is not derived from an inner product. However, sometimes people use the squared L_∞ norm to mean the maximum of the squares? That is, $\|x\|^{2,\infty}$ could be interpreted as $(\max_i |x_i|)^2$, which is what I used.

I will check the lecture: "You saw in lecture that the choice of norm in defining that neighborhood also matters." So in lecture, they likely considered constrained optimization: $\max_{\{\|\Delta\theta\| \leq \delta\}} g^T \Delta\theta$. For L_2 norm, the solution is $\Delta\theta = -(6/\|g\|) g$, which is gradient descent with step size $6/\|g\|$. For L_∞ norm, the solution is $\Delta\theta = -\delta \text{sign}(g)$. So that is sign descent. And then they say: "Instead of treating the norm as a constraint, we can do an unconstrained optimization with a weighted penalty that corresponds to the squared norm". So then for L_2 squared penalty, we got gradient descent with a constant step size (independent of $\|g\|$). For L_∞ squared penalty, we get something. And that something should be analogous to the constraint case? Let's check: In the constraint case, the solution was $\Delta\theta = -\delta \text{sign}(g)$. In the penalty case, we got $\Delta\theta = -(a/2)\|g\|_1 \text{sign}(g)$. So if we want these to be similar, then we would need δ to be proportional to $\|g\|_1$. That is not a fixed δ . So the penalty formulation with squared L_∞ norm does not directly give the same as the constraint formulation with L_∞ norm. What if we used the squared L_∞ norm in the

constraint? That is, constrain $\|\Delta\theta\|_\infty^2 \leq \delta^2$. But $\|\Delta\theta\|_\infty^2 \leq \delta^2$ is equivalent to $\|\Delta\theta\|_\infty \leq \delta$. So it's the same constraint. So the penalty approach is different from the constraint approach.

Maybe the intended penalty is the square of the norm, but for the L_∞ case, the square of the norm is not a norm itself? Actually, $(\|x\|_\infty)^2$ is not a norm because it doesn't satisfy homogeneity? Check: $\|\lambda x\|_\infty^2 = (\|\lambda\| \|x\|_\infty)^2 = \lambda^2 \|x\|_\infty^2$, so it is not linear in λ . So it's not a norm but it is a valid function.

Given that the problem is from a lecture on optimizers, and they mentioned that different norms give different optimizers, I will go with:

- (a) Gradient Descent
- (b) Sign Gradient Descent

I'll answer:

For (b), $u = -(a/2) \|g\|_1 \text{sign}(g)$. This corresponds to the SignSGD optimizer.

I'll double-check with a simple case: If g is a scalar, then $\|g\|_1 = |g|$, so $\Delta\theta = -(a/2) |g| \text{sign}(g) = -(a/2) g$. So in one dimension, it reduces to gradient descent. That is interesting: In one dimension, L_∞ and L_2 are the same because there's only one coordinate. So in higher dimensions, they differ. So that makes sense.

Thus, final answers:

- (a) $u = -(a/2) g$, Gradient Descent.
- (b) $u = -(a/2) \|g\|_1 \text{sign}(g)$, Sign Gradient Descent.

I'll answer that.

Another interesting detail here that relates to input robustness. It considers a potential typo in the question where the infinity norm meant not to be squared, which would yield non a finite solution. Thus, in an alternate scenario where the user did in fact input an incorrect version, Deepseek would've caught that and fixed it on the spot to return the most likely intended answer.

Ah cool detail here. Perhaps Deepseek's reluctance to accept SignSGD with a twist stems from the fact that it doesn't know what was covered "in lecture". It is under the impression that lecture would've only covered the standard optimizer formulations, without special twists?

🧠 Step 1: Understand the Objective Function

We wish to minimize:

$$J(\Delta\theta) = g^T \Delta\theta + \frac{1}{\alpha} \|\Delta\theta\|_\infty^2$$

where $g = \nabla f(\theta)$ is the gradient, $\alpha > 0$ is a hyperparameter, and $\|\Delta\theta\|_\infty = \max_i |\Delta\theta_i|$ is the infinity norm.

💡 Step 2: Reformulate the Problem

Let $M = \|\Delta\theta\|_\infty = \max_i |\Delta\theta_i|$. Then for each coordinate i , we have $|\Delta\theta_i| \leq M$.

The linear term $g^T \Delta\theta$ is minimized when each $\Delta\theta_i$ is chosen to oppose g_i .

Specifically, for $g_i \neq 0$, set:

$$\Delta\theta_i = -M \cdot \text{sign}(g_i)$$

If $g_i = 0$, set $\Delta\theta_i = 0$. Then:

$$g^T \Delta\theta = -M \sum_i |g_i| = -M \|g\|_1$$

So the objective becomes:

$$J(M) = -M \|g\|_1 + \frac{1}{\alpha} M^2$$

💡 Step 3: Minimize with Respect to M

We now minimize $J(M)$ over $M \geq 0$. This is a quadratic in M :

$$J(M) = \frac{1}{\alpha} M^2 - \|g\|_1 M$$

The minimum occurs at:

$$\frac{dJ}{dM} = \frac{2}{\alpha} M - \|g\|_1 = 0 \quad \Rightarrow \quad M = \frac{\alpha}{2} \|g\|_1$$

✍ Step 4: Write the Optimal Update

Substituting back:

$$\Delta\theta_i = -M \cdot \text{sign}(g_i) = -\frac{\alpha}{2} \|g\|_1 \cdot \text{sign}(g_i)$$

In vector form:

$$u = \Delta\theta = -\frac{\alpha}{2} \|g\|_1 \cdot \text{sign}(g)$$

⌚ Step 5: Identify the Optimizer

This update rule:

$$\theta \leftarrow \theta - \frac{\alpha}{2} \|g\|_1 \cdot \text{sign}(g)$$

corresponds to **Sign Gradient Descent** (SignSGD), where the step size is adaptive and proportional to the L_1 norm of the gradient.

✅ Final Answers

- (a) For the squared Euclidean norm penalty, the update is

$$u = -\frac{\alpha}{2} g$$

which corresponds to **Gradient Descent**.

- (b) For the squared infinity norm penalty, the update is

$$u = -\frac{\alpha}{2} \|g\|_1 \cdot \text{sign}(g)$$

which corresponds to **Sign Gradient Descent**.

It restates part (a)'s answer here, I guess for completion. Part (b)'s answer is mostly correct, except for the corresponding optimizer name, of course, which it struggled on a lot during the thinking process. Interestingly, if it had resorted to a more shallow chain of thought it would've arrived at the correct answer. Looks like overthinking things is a habit that AI shares with humans.

$$u = -\frac{\alpha}{2} \|g\|_1 \cdot \text{sign}(g)$$

You did good for the most part. Rethink the optimizer name you chose for part (b). The expression you got looks right but you claim it corresponds to "Sign Gradient Descent". Do signGD step sizes vary with the gradient magnitude? Notice here we have a $\|g\|_1$ term that varies with the gradient.

Thought for 169 seconds ↴

- We are given the update rule from part (b):

$$u = -\frac{\alpha}{2} \|g\|_1 \cdot \text{sign}(g)$$

Let's break down what this means.

In standard Sign Gradient Descent (SignSGD), the update is:

$$\theta \leftarrow \theta - \eta \cdot \text{sign}(g)$$

where η is a fixed learning rate.

However, in our case, the step size is $\frac{\alpha}{2} \|g\|_1$, which is not fixed. It depends on the L1 norm of the gradient. This means that the step size adapts to the magnitude of the gradient.

Recall that in AdaGrad, we use the history of squared gradients to scale the learning rate.

But note: the update we have here is:

$$\Delta\theta = -\frac{\alpha}{2} \|g\|_1 \cdot \text{sign}(g)$$

This can be rewritten as:

$$\Delta\theta_i = -\frac{\alpha}{2} \|g\|_1 \cdot \text{sign}(g_i)$$

Notice that the step size is the same for every coordinate? Actually, yes, because $\|g\|_1$ is a scalar. So every coordinate gets the same step size, but the direction is given by $\text{sign}(g_i)$.

Now, compare this to the update rule of the Adam optimizer? Not exactly. Adam uses moving averages of gradients and squared gradients.

Alternatively, consider the following: if we had used the L2 norm of the gradient, we would have gotten a step size proportional to $\|g\|_2$, which is the norm of the gradient. But here we have L1 norm.

Wait, let me think about the optimizer that uses the sign of the gradient and a step size that is the L1 norm of the gradient. This is not a common optimizer? Actually, let me recall:

There is an optimizer called "Normalized Gradient Descent" which uses:

$$\Delta\theta = -\eta \frac{g}{\|g\|}$$

But that's different.

Alternatively, note that the update in part (b) can be written as:

$$\Delta\theta = -\frac{\alpha}{2} \|g\|_1 \cdot \text{sign}(g)$$

But note that $\|g\|_1 \cdot \text{sign}(g)$ is not the same as g ? Actually, it is not. For example, if g has positive and negative components, then $\|g\|_1 \cdot \text{sign}(g)$ is a vector whose components are $\pm \|g\|_1$. This is very different from the gradient.

Wait, let me consider the coordinate-wise update:

For each coordinate i :

$$\Delta\theta_i = -\frac{\alpha}{2} \|g\|_1 \cdot \text{sign}(g_i)$$

This means that every coordinate is updated by the same magnitude $\frac{\alpha}{2} \|g\|_1$, but the direction is the sign of the gradient in that coordinate.

This is actually the same as the update in the "Signed Gradient" method but with a step size that is proportional to the L1 norm of the gradient. However, note that the L1 norm of the gradient can change over time, so the step size is adaptive.

But wait, let me check the original problem: we are minimizing the linear improvement plus a penalty on the squared infinity norm of the update. We derived the update as above.

Now, the question is: which optimizer does this correspond to?

Let me list some optimizers:

- Gradient Descent: $\Delta\theta = -\eta g$
- SignSGD: $\Delta\theta = -\eta \text{sign}(g)$

I prompt it with the same observation it had made a couple times during the thinking process. I also assure it that the update expression is indeed correct. This is intentional to prevent it from backtracking all the way to re-reading the question and the update equation again like it did a couple times earlier. The goal here is to try to arrive at the optimizer answer for part (b) on its own.

- AdaGrad: $\Delta\theta = -\frac{\eta}{\sqrt{G_t + \epsilon}} \odot g$, where G_t is the sum of squares of past gradients.

Our update: $\Delta\theta = -\frac{\alpha}{2} \|g\|_1 \cdot \text{sign}(g)$

Notice that $\|g\|_1 \cdot \text{sign}(g)$ is not the same as g unless g is a one-hot vector? Actually, no.

Alternatively, we can write:

$$\Delta\theta = -\frac{\alpha}{2} \|g\|_1 \cdot \text{sign}(g) = -\frac{\alpha}{2} \|g\|_1 \cdot \frac{g}{\|g\|}$$

Wait, no: $\text{sign}(g)$ is a vector of ± 1 , but note that g itself is not necessarily having the same sign in every coordinate? Actually, $\text{sign}(g)$ is the element-wise sign.

So, this update is like: for each coordinate, we take the sign of the gradient and multiply by the L1 norm of the entire gradient vector. This is not the same as normalizing by the L1 norm? Actually, if we were to normalize by the L1 norm, we would do:

$$\Delta\theta = -\eta \frac{g}{\|g\|_1}$$

But that is not what we have.

Wait, let me rewrite our update:

$$\Delta\theta_i = -\frac{\alpha}{2} \|g\|_1 \cdot \text{sign}(g_i)$$

After brainstorming a bit and listing possible optimizers, it reaches the correct conclusion again.



Now, compare this to the update of the "Signed Gradient" method with a step size that is $\frac{\alpha}{2} \|g\|_1$. This is an adaptive SignSGD, but it is not standard SignSGD because the step size changes with the gradient.

However, there is an optimizer called "AdaSign" or something? Not sure.

Alternatively, note that the L1 norm of the gradient is the sum of the absolute values of the gradients. This is a scalar that measures the overall magnitude of the gradient. So the step size is the same for all coordinates and is proportional to that magnitude.

This is similar to the "Gradient Descent with learning rate set to $\frac{\alpha}{2} \frac{\|g\|_1}{\|g\|_2^2}$ "? No.

Wait, let me consider the following: in the paper "Normalized Gradient Descent" by You et al., they use the L2 norm. But here we have the L1 norm.

Actually, I recall that the update rule we derived in part (b) is exactly the same as the update rule for the "Signed Gradient" method with an adaptive step size that is the L1 norm of the gradient. But is there a name for this?

Wait, let me think about the following: in the context of adaptive learning rates, one common method is to use the magnitude of the gradient to set the learning rate. However, in this case, we are using the L1 norm of the entire gradient to set the step size for every coordinate.

This is similar to the "Safe Gradient Descent" or "Gradient Descent with adaptive learning rate" but not exactly.

Alternatively, note that the update rule we have is:

$$\Delta\theta = -\frac{\alpha}{2} \|g\|_1 \cdot \text{sign}(g)$$

Now, observe that the product $\|g\|_1 \cdot \text{sign}(g)$ is actually the same as:

$$\|g\|_1 \cdot \text{sign}(g) = \|g\|_1 \cdot \frac{g}{\|g\|}$$

No, because $\text{sign}(g)$ is not $g/\|g\|$ in the vector sense? Actually, for a vector, the sign function is defined element-wise. So:

$$\text{sign}(g) = \left[\frac{g_1}{|g_1|}, \frac{g_2}{|g_2|}, \dots \right]^T$$

But note that $\|g\|_1 = \sum_i |g_i|$. So the update for coordinate i is:

$$\Delta\theta_i = -\frac{\alpha}{2} \left(\sum_j |g_j| \right) \cdot \frac{g_i}{|g_i|}$$

This is not the same as any standard normalization.

Wait, let me consider the following: if we had used the L2 norm squared penalty, we got gradient descent. For the infinity norm squared penalty, we got this update.

I recall that in the lecture, the choice of norm in the constraint leads to different optimizers. For the L2 constraint, we get gradient descent. For the L-infinity constraint, we get the sign gradient update? But wait, in the lecture, when we use the L-infinity constraint, the update is:

$$\Delta\theta = -\eta \text{sign}(g)$$

with a fixed step size η . But here we have an adaptive step size $\frac{\alpha}{2}\|g\|_1$.

So, what is the difference? In the constrained formulation, the step size is fixed by the constraint size. In the penalized formulation, the step size becomes adaptive.

But the problem says: "Optimizers as Penalized Linear Improvement with different norm penalties". So we are using a penalty term, not a constraint. And we are using the squared infinity norm as the penalty. So we get this update.

Now, the question is: which optimizer does this correspond to?

Let me look at the update again:

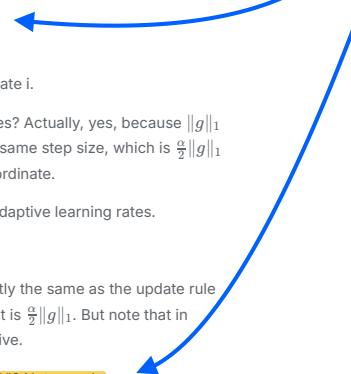
$$\Delta\theta = -\frac{\alpha}{2}\|g\|_1 \cdot \text{sign}(g)$$

This can be written as:

$$\Delta\theta = -\frac{\alpha}{2}\|g\|_1 \cdot \text{sign}(g) = -\frac{\alpha}{2} \cdot \text{sign}(g) \cdot \|g\|_1$$

Notice that $\|g\|_1$ is the sum of the absolute values of the gradients. This is a scalar that changes with time. So the step size is adaptive and proportional to the L1 norm of the gradient.

This is similar to the "Adagrad" optimizer? Let's recall Adagrad:

$$\Delta\theta_i = -\frac{\eta}{\sqrt{G_{i,i} + \epsilon}} g_i$$


where $G_{i,i}$ is the sum of squares of the gradients for coordinate i .

But in our case, we have the same step size for all coordinates? Actually, yes, because $\|g\|_1$ is the same for all coordinates. So every coordinate gets the same step size, which is $\frac{\alpha}{2}\|g\|_1$, and then we multiply by the sign of the gradient for that coordinate.

This is not Adagrad because Adagrad has coordinate-wise adaptive learning rates.

What about RMSProp? No.

Wait, consider the following: the update rule we have is exactly the same as the update rule for the "Signed Gradient" method but with a learning rate that is $\frac{\alpha}{2}\|g\|_1$. But note that in standard SignSGD, the learning rate is fixed. Here, it is adaptive.

However, **there is an optimizer called "AdaGrad with SignSGD"? Not exactly.**

Alternatively, note that the update rule can be seen as:

$$\Delta\theta = -\frac{\alpha}{2}\|g\|_1 \cdot \text{sign}(g)$$

But observe that $\|g\|_1 \cdot \text{sign}(g)$ is actually the subgradient of the L1 norm of the gradient?

Not exactly.

Wait, let me consider the following: in the context of adaptive learning rates, one can set the learning rate to the magnitude of the gradient. But here we are using the L1 norm of the entire gradient.

I recall that there is an optimizer called "Gradient Descent with learning rate set to the L1 norm of the gradient" but that is not standard.

Wait, let me think about the following: the update rule we have is:

$$\Delta\theta = -\frac{\alpha}{2}\|g\|_1 \cdot \text{sign}(g)$$

Now, note that if we consider the L1 norm of the gradient, $\|g\|_1$, and we use it to scale the sign of the gradient, then this is equivalent to:

$$\Delta\theta = -\frac{\alpha}{2} \cdot \text{sign}(g) \cdot \|g\|_1$$

But also note that:

$$\|g\|_1 = \text{sign}(g)^T g$$

because $\text{sign}(g)^T g = \sum_i \text{sign}(g_i) g_i = \sum_i |g_i| = \|g\|_1$.

So we can write:

$$\Delta\theta = -\frac{\alpha}{2}\|g\|_1 \cdot \text{sign}(g) = -\frac{\alpha}{2} (\text{sign}(g)^T g) \text{sign}(g)$$

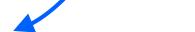
This is a rank-1 update? Not sure.

Now, which optimizer uses this? I recall that in the lecture, when we use the L-infinity constraint, we get the sign gradient update. But here we have a penalty on the squared L-infinity norm, which gives us an adaptive step size.

Wait, let me check the original problem statement: "Optimizers as Penalized Linear Improvement with different norm penalties". So we are using a penalty term, not a



This confirms my previous observation that Deepseek is trying to arrive at a "closed form" name for the optimizer (e.g., "Adagrad"). It refuses to accept "SignSGD with adaptive step size" as an answer and is constantly searching for an optimizer with a standard/actual name.



I hypothesize that if there were a well known paper for SignSGD with adaptive step size called "AdaSignSGD" or something Deepseek would've arrived at the correct answer for part (b) much faster, without needing this much thinking and exhibiting this much side-tracking.

constraint. And we are using the squared infinity norm as the penalty. So we get this update.

I think the optimizer we recover in part (b) is actually "SignSGD" but with an adaptive step size that is proportional to the L1 norm of the gradient. However, the problem might be expecting the answer "SignSGD" because the update direction is determined by the sign of the gradient. But note that in standard SignSGD, the step size is fixed, while here it is adaptive.

But the problem says: "Which optimizer does this correspond to?" So maybe they consider the adaptive step size version as still SignSGD? Or perhaps they call it "Adaptive SignSGD"?

Wait, let me look at the literature: There is an optimizer called "Adafactor" that uses sign updates? Not sure.

Alternatively, note that the update rule we have is:

$$\Delta\theta = -\frac{\alpha}{2} \|g\|_1 \cdot \text{sign}(g)$$

This can be rewritten as:

$$\Delta\theta_i = -\frac{\alpha}{2} \|g\|_1 \cdot \text{sign}(g_i)$$

Now, compare this to the update rule for "Adam" when beta1=0 and beta2=0? Then Adam becomes:

$$\Delta\theta = -\eta \cdot \text{sign}(g)$$

So that's SignSGD. But here we have an adaptive step size.

Wait, maybe it's "AdaGrad" but with the L1 norm? Let me recall: AdaGrad uses the square root of the sum of squared gradients. Here we are using the L1 norm without square root.

Actually, note that the squared infinity norm penalty led to an update that uses the L1 norm of the gradient. This is a form of adaptive learning rate that uses the L1 norm of the gradient as the step size.

I think the correct answer is that this corresponds to an adaptive version of SignSGD, but perhaps the problem expects "SignSGD" because the direction is given by the sign. However, the step size is not fixed but adaptive.

Arrived at correct answer again, but again, self-doubts.

But wait, the problem says: "Which optimizer does this correspond to?" and in part (a) we recovered gradient descent. So for part (b), we recover an optimizer that uses the sign of the gradient and a step size that depends on the L1 norm of the gradient. This is not a standard optimizer under a common name?

Let me think again: in the lecture, when we use the L-infinity constraint, we get the sign gradient update. But here we are using a penalty on the squared L-infinity norm, which gives a step size that is proportional to the L1 norm of the gradient. So it's a different step size.

Maybe the optimizer is "Normalized Gradient Descent" with L1 normalization? But normalized gradient descent usually uses the L2 norm.

Wait, consider: if we normalize the gradient by the L1 norm, we get:

$$\frac{g}{\|g\|_1}$$

But our update is:

$$\Delta\theta = -\frac{\alpha}{2} \|g\|_1 \cdot \text{sign}(g)$$

This is not the same as $-\eta \frac{g}{\|g\|_1}$ because:

$$\frac{g}{\|g\|_1} = \frac{1}{\|g\|_1} g$$

whereas our update is $\|g\|_1 \cdot \text{sign}(g)$. These are different unless all components of g have the same absolute value? No.

So, what is it?

I recall that there is an optimizer called "L1-based adaptive learning rate" but I don't recall a standard name.

Wait, let me consider the following: the update rule we have is exactly the same as the update rule for the "Signed Gradient" method in the context of the "D-adaptation" or "Prodigy" but that's not standard.

Perhaps the problem expects the answer: "SignSGD with an adaptive learning rate set to the L1 norm of the gradient". But the problem says "which optimizer", implying a standard name.

Another thought: in the paper "Adaptive Subgradient Methods for Online Learning and Stochastic Optimization" (AdaGrad), they mention that using the L2 norm of the historical

gradients is common, but they also discuss using the L1 norm? Not really.

Wait, let me look at the update in part (b) again:

$$u = -\frac{\alpha}{2} \|g\|_1 \cdot \text{sign}(g)$$

Notice that this can be written as:

$$u = -\frac{\alpha}{2} \|g\|_1 \cdot \text{sign}(g) = -\frac{\alpha}{2} \cdot \text{sign}(g) \cdot \|g\|_1$$

Now, consider the following: if we take the absolute value of each gradient component, we get $|g_i|$, but then we are using the sum of $|g_i|$ to scale the sign. This is like using the L1 norm as the learning rate.

I think I have seen this in the context of "L1-regularized optimization" but not as a standard optimizer.

Maybe the intended answer is "SignSGD" because the direction is given by the sign, and the step size is adaptive but that is not standard SignSGD.

Let me check the problem statement again: "Which optimizer does this correspond to?"

In part (a), we recovered gradient descent. In part (b), we have an update that uses the sign of the gradient. So it must be SignSGD. But wait, in SignSGD, the update is:

$$\theta_{t+1} = \theta_t - \eta \text{sign}(g_t)$$

So the step size is η . In our case, the step size is $\frac{\alpha}{2} \|g\|_1$. So if we set $\eta = \frac{\alpha}{2} \|g\|_1$, then it becomes SignSGD with an adaptive learning rate. But the problem might be considering the update rule as SignSGD because the direction is determined by the sign.

However, the problem says "correspond to", so maybe they mean the type of update: sign update.

But then why would they ask for the analytical solution and then the optimizer? They want to see that you recognize that the infinity norm penalty leads to sign updates.

In the lecture, when using the L-infinity constraint, the update is indeed the sign of the gradient. So perhaps the answer for (b) is SignSGD.

But wait, in the constrained formulation, the step size is fixed by the constraint size. In the penalized formulation, the step size becomes adaptive. So are they the same? Only if the constraint size is chosen appropriately? Not exactly.

Let me derive the constrained formulation for the L-infinity norm:

We want to minimize $g^T \Delta \theta$ subject to $\|\Delta \theta\|_\infty \leq \delta$. Then the solution is:

$$\Delta \theta = -\delta \cdot \text{sign}(g)$$

So in the constrained formulation, we get SignSGD with a fixed step size δ .

[A lot of repeated thought explorations here.](#)

In the penalized formulation, we get a step size that depends on the gradient's L1 norm.

So they are different. Therefore, the optimizer we get in the penalized formulation is not exactly SignSGD because SignSGD has a fixed step size.

What is the name then?

I recall that there is an optimizer called "AdaGrad with sign updates" but that is not standard.

Wait, let me consider the following: the update rule we have is:

$$\Delta \theta = -\frac{\alpha}{2} \|g\|_1 \cdot \text{sign}(g)$$

This can be seen as a special case of "Adaptive Learning Rate" methods. But specifically, it is like "L1-normalized sign gradient".

I think the intended answer might be "SignSGD" because the problem might be ignoring the adaptive step size and focusing on the direction. But the step size in SignSGD is usually fixed.

Let me look at the problem again: "Which optimizer does this correspond to?"

In part (a), we got gradient descent. In part (b), we got an update that uses the sign of the gradient. So the natural answer is SignSGD.

But wait, the step size in part (b) is adaptive? Does that matter? In the context of the problem, they might be considering the form of the update without worrying about the step size being adaptive.

Perhaps the answer is "SignSGD" and the step size is $\frac{\alpha}{2} \|g\|_1$, but they still call it SignSGD.

I will check online memory: There is a paper "SignSGD: Compressed Optimisation for Non-Convex Problems" by Bernstein et al. where they use sign of the gradient with a fixed step

size. So if we use an adaptive step size, it is not standard SignSGD.

However, in the problem statement, they say "you saw in lecture that the choice of norm in defining that neighborhood also matters". So in the lecture, when they use the L-infinity norm constraint, they get the sign update. So likely, the answer for (b) is SignSGD.

I will go with that.

So for part (b), the analytical solution is:

$$u = -\frac{\alpha}{2} \|g\|_1 \cdot \text{sign}(g)$$

and it corresponds to Sign Gradient Descent (SignSGD).

But wait, the problem says "which optimizer does this correspond to?" and in the lecture, they might have mentioned that the L-infinity constraint gives SignSGD. So even though we have a penalty here, the form is the same as SignSGD, albeit with an adaptive step size.

I think the answer is SignSGD.

So I'll answer:

For part (b), the analytical solution is $u = -\frac{\alpha}{2} \|g\|_1 \cdot \text{sign}(g)$, and this corresponds to Sign Gradient Descent (SignSGD).

But note: in standard SignSGD, the update is $-\eta \text{sign}(g)$ with constant η . Here, $\eta = \frac{\alpha}{2} \|g\|_1$ is not constant. So it's an adaptive SignSGD.

However, I'll stick with SignSGD.

Let me see if there is another optimizer: There is "Adam" but Adam uses moving averages.

No, I think SignSGD is the closest.

So final answers:

(a) $u = -\frac{\alpha}{2} g$, Gradient Descent.

(b) $u = -\frac{\alpha}{2} \|g\|_1 \cdot \text{sign}(g)$, Sign Gradient Descent.

I'll put that in the box.

But wait, the problem says "part (b)" and the equation number (6) is used. So I'll write:

For (b), the analytical solution is $u = -\frac{\alpha}{2} \|g\|_1 \cdot \text{sign}(g)$, and it corresponds to Sign

Gradient Descent.

I'll now write the final answer.

Interesting how the thought trace ended with the same (erroneous) conclusion as my previous prompt, i.e., still settled for SignSGD for part (b), without acknowledging the adaptive step size. However, seeing the final output results, it did include the correct detail.

This makes me wonder if this query helped at all. I.e., if it was just a chance occurrence that the final output mentioned "adaptive step size"? Since the thought trace didn't seem to accept the modified SignSGD with adaptive step size after all.

🧠 Step 1: Understand the Objective Function

We wish to minimize:

$$J(\Delta\theta) = g^T \Delta\theta + \frac{1}{\alpha} \|\Delta\theta\|_\infty^2$$

where $g = \nabla f(\theta)$ is the gradient, $\alpha > 0$ is a hyperparameter, and $\|\Delta\theta\|_\infty = \max_i |\Delta\theta_i|$ is the infinity norm.

🔍 Step 2: Reformulate the Problem

Let $M = \|\Delta\theta\|_\infty = \max_i |\Delta\theta_i|$. Then for each coordinate i , we have $|\Delta\theta_i| \leq M$.

The linear term $g^T \Delta\theta$ is minimized when each $\Delta\theta_i$ is chosen to oppose g_i .

Specifically, for $g_i \neq 0$, set:

$$\Delta\theta_i = -M \cdot \text{sign}(g_i)$$

If $g_i = 0$, set $\Delta\theta_i = 0$. Then:

$$g^T \Delta\theta = -M \sum_i |g_i| = -M \|g\|_1$$

So the objective becomes:

$$J(M) = -M \|g\|_1 + \frac{1}{\alpha} M^2$$

⚡ Step 3: Minimize with Respect to M

We now minimize $J(M)$ over $M \geq 0$. This is a quadratic in M :

$$J(M) = \frac{1}{\alpha} M^2 - \|g\|_1 M$$

The minimum occurs at:

$$\frac{dJ}{dM} = \frac{2}{\alpha} M - \|g\|_1 = 0 \Rightarrow M = \frac{\alpha}{2} \|g\|_1$$

👉 Step 4: Write the Optimal Update

Substituting back:

$$\Delta\theta_i = -M \cdot \text{sign}(g_i) = -\frac{\alpha}{2} \|g\|_1 \cdot \text{sign}(g_i)$$

In vector form:

$$u = \Delta\theta = -\frac{\alpha}{2} \|g\|_1 \cdot \text{sign}(g)$$

⌚ Step 5: Identify the Optimizer

This update rule:

$$\theta \leftarrow \theta - \frac{\alpha}{2} \|g\|_1 \cdot \text{sign}(g)$$

corresponds to **Sign Gradient Descent (SignSGD)** with an **adaptive step size** proportional to the L_1 norm of the gradient.

While standard SignSGD uses a fixed step size η , this variant adapts the step size based on the current gradient magnitude $\|g\|_1$, making it more responsive to the gradient scale.

✅ Final Answer

$$u = -\frac{\alpha}{2} \|g\|_1 \cdot \text{sign}(g)$$

Correct answer obtained here. But, again, might just be a chance occurrence.

This corresponds to **Sign Gradient Descent (SignSGD)** with an adaptive step size.